

# MITIGATING SELECTION BIAS WITH NODE PRUNING AND AUXILIARY OPTIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models (LLMs) often show unwarranted preference for certain choice options when responding to multiple-choice questions, posing significant reliability concerns in LLM-automated systems. To mitigate this selection bias problem, previous solutions utilized debiasing methods to adjust the model’s input and/or output. Our work, in contrast, investigates the model’s internal representation of the selection bias. Specifically, we introduce a novel debiasing approach, Bias Node Pruning (BNP), which eliminates the linear layer parameters that contribute to the bias. Furthermore, we present Auxiliary Option Injection (AOI), a simple yet effective input modification technique for debiasing, which is compatible even with black-box LLMs. To provide a more systematic evaluation of selection bias, we review existing metrics and introduce Choice Kullback-Leibler Divergence (CKLD), which addresses the insensitivity of the commonly used metrics to imbalance in choice labels. Experiments show that our methods are robust and adaptable across various datasets when applied to three LLMs.

## 1 INTRODUCTION

The advent of large language models (LLMs) has revolutionized artificial intelligence applications, particularly in the domain of natural language processing. These models have demonstrated outstanding performance across a variety of use cases, including chatbots, machine translation, text generation, data annotation, etc. Their ability to answer questions with high precision has opened up new avenues for automated systems.

Despite their remarkable abilities, LLMs suffer from the selection bias problem that often occurs in answering multiple-choice questions (MCQs). When selecting the answer for an MCQ, many LLMs prefer the choices in a given position (*e.g.*, the last choice), or with a specific choice symbol (*e.g.*, (A) or (3)) (Zheng et al., 2024; Wei et al., 2024; Pezeshkpour & Hruschka, 2024). This phenomenon degrades model performance.

Many previous works have attempted to explain this phenomenon and/or propose diverse ways to mitigate selection bias. While there are a few works focused on either modifying the input format (Li et al., 2023b; Robinson et al., 2023) or calibrating the output probabilities (Zheng et al., 2024; Reif & Schwartz, 2024; Wei et al., 2024), to the best of our knowledge, no embedding or parameter-level investigation has been performed. Because selection bias originates from internal parameter-level computations, it is crucial to explore how the LLM embeddings contribute to the bias in their output responses.

Understanding the internal representation of selection bias can help us combat it. By scrutinizing the interaction between the internal representation and the LLM parameters, we develop a novel approach to debias the model. Specifically, we propose **Bias Node Pruning** (BNP), which eliminates nodes in the final linear layer that contribute to selection bias. By dropping as few as 32 out of 4096

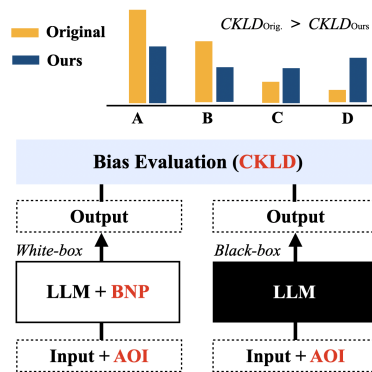


Figure 1: We propose BNP and AOI to reduce selection bias for white-box and black-box models. The CKLD metric is also proposed to encourage a more standardized evaluation of the bias.

nodes in the final layer, we can significantly reduce selection bias and improve question-answering performance. In addition, we find that introducing an “I don’t know” option in the input reduces bias and enhances task performance. This **Auxiliary Option Injection** (AOI) technique is a simple method that can be applied to even black-box scenarios.

Although mitigating selection bias is an important task, even quantifying the extent of selection bias is in itself a difficult problem. Previous research has adopted several bias evaluation metrics, such as the Standard Deviation of Recalls (RStd) (Zheng et al., 2024) and the Relative Standard Deviation (RSD) (Reif & Schwartz, 2024). However, these metrics are insensitive to imbalance of choices, which can lead them to incorrectly indicate selection bias when none exists. To address this concern, we propose the **Choice Kullback-Leibler Divergence (CKLD)**, a novel bias evaluation metric that is sensitive to the imbalance. Figure 1 depicts our contributions to the overall pipeline.

We conducted experiments and analyses to evaluate the debiasing performance of our methods, adopting the proposed CKLD metric. We validate the efficacy of our approach on widely used public benchmark datasets with various LLMs. Results show that our method generally improves debiasing and task performance, and can be utilized together with other baseline methods (*e.g.*, Chain-of-Thought, In-Context Learning, or Decoding by Contrasting Layers).

Our **contributions** are four-fold. In this work, we:

- Propose *Bias Node Pruning* (BNP), a novel debiasing approach that removes parameters from the final linear layer that contribute to selection bias.
- Introduce *Auxiliary Option Injection* (AOI), which is a simple prompting tactic for MCQ answering. Along with BNP, our debiasing methods improve accuracy by upto 24.9%.
- Review existing metrics to systematically evaluate selection bias, and introduce *Choice Kullback-Leibler Divergence* (CKLD) to address their weakness with imbalanced labels.
- Underscore the broad applicability of our approach to various baselines and also demonstrate that our AOI method can debias black-box large language models.

## 2 SELECTION BIAS IN LLMs

Although LLMs are most often used for text generation, some tasks involve responding to multiple-choice questions (MCQs). For example, LLMs are increasingly used to annotate data samples, a task that requires selecting the best choice from several options. When responding to MCQs, however, LLMs suffer from selection bias, which is the model’s inclination to prefer a choice option bound with a specific symbol or located in a certain order. In this section, we formally define selection bias (§ 2.1) and discuss when and where the signs of selection bias are observed (§ 2.2).

### 2.1 SELECTION BIAS PROBLEM

Selection bias refers to a model’s tendency to select options in a given position or with a given symbol among MCQ choices, regardless of the correctness of its choice. This includes the model’s *a priori* preference of a certain choice symbol, and its inclination to favor the choice presented at a specific ordering position (Zheng et al., 2024). **In this work, we define selection bias as the discrepancy between the model’s selection for the original choice ordering of a question and its expected option selection across all possible choice orderings of a question.** If the model consistently selects the same choice option (*i.e.*, the content of the option) regardless of its position, the discrepancy is zero, indicating no selection bias. Conversely, a high level of selection bias suggests that the model’s selection for certain choice orderings may deviate from expectation.

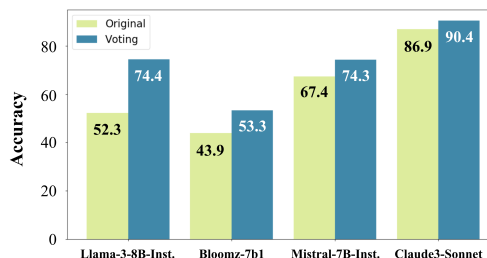


Figure 2: Comparison of the original and voting accuracy with different LLMs via zero-shot querying. Note, Claude3-Sonnet is evaluated under the black-box setting (Section 5.2)

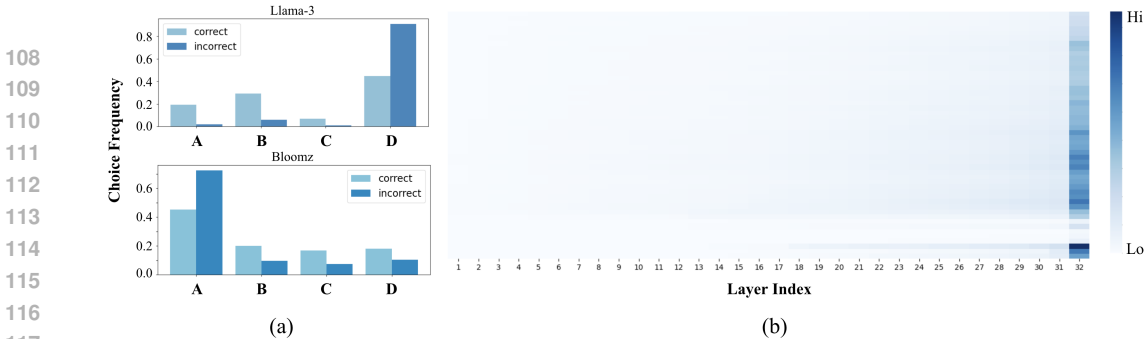


Figure 3: (a) Choice frequency tends to have a sharper distribution when the model’s response is incorrect. (b) In Llama-3, selection bias is predominantly observed to be in the final output layer of the decoder. Other model figures are in Appendix D.

**Empirical demonstration.** Motivated by this definition, Figure 2 shows the existence of selection bias on four LLMs. The lighter bars show each model’s accuracy on the ARC-Challenge dataset (Clark et al., 2018). The darker bars, on the other hand, show the accuracy of the answers retrieved by majority voting across all possible choice permutations, which can be interpreted as the expected output across all choice orderings. If the model is free of selection bias, voting will always output the same choice as the original question, rendering the same accuracy in all cases. If the model entails selection bias, on the other hand, its response to the original question may deviate from the expected response, leading to a bigger gap between the voting accuracy and original accuracy. In the figure, selection bias exists with all four models and is greatest with Llama-3.

## 2.2 MOTIVATING ANALYSES

While selection bias is a prevalent problem in querying the large language model (LLM), it is important to properly identify when and where the bias is captured. Here, we provide two simple analyses that motivate the design of our debiasing methods.

**Selection bias is prominently captured when the model is incorrect.** Figure 3(a) shows the frequency of choices of the four response options on the ARC-Challenge dataset (Clark et al., 2018) using Llama-3-8B-Instruct (Meta, 2024) and Bloomz-7b1 (Muennighoff et al., 2023). We manipulated the test dataset to include all possible orderings of the MCQ choices. Thus, the bars should be at 0.25. However, the models prefer answer choices ‘D’ and ‘A’, respectively. These preferences are pronounced in cases where the models produce incorrect responses, as opposed to correct ones. This observation highlights the role of selection bias in incorrect predictions and motivates our focus on analyzing cases where the model’s output is incorrect.

**Selection bias is prominently observed in the final decoder layers.** To capture the selection bias, we investigate the difference between the correct and incorrect sample embeddings extracted from different locations. Specifically, we explore the discrepancies within a single sample by permuting the sequence of choices in the question. The difference between the embeddings within the choice-permuted set removes the sample-specific semantic information while the pure effect of the selection bias remains in the difference.

Accordingly, we first retrieve the intermediate embeddings of an LLM by computing the  $t$ -th token embedding from the  $\ell$ -th decoder layer as  $\mathbf{z}_{\ell,t} = f_{\ell}(\mathbf{x}_{\mathcal{A}})_t$ , where  $f_{\ell}$  is the LLM decoder up to the  $\ell$ -th layer and  $\mathbf{x}_{\mathcal{A}}$  is the input with answer choices  $\mathcal{A}$ . For brevity of notation, let  $\mathbf{z} \in \mathbb{R}^d$  be the embedding from an arbitrary layer and token location. Then, we quantify the selection bias by computing the embedding difference between the correct and incorrect questions within the permutations of  $\mathcal{A}$ . That is, the bias vector  $\mathbf{b}$  for a sample  $\mathbf{x}$  is defined as

$$\mathbf{b}_{\mathbf{x}} = \frac{1}{n_-} \sum_{i=1}^{n_-} \mathbf{z}_-^{(i)} - \frac{1}{n_+} \sum_{i=1}^{n_+} \mathbf{z}_+^{(i)}, \tag{1}$$

where  $\mathbf{z}_-$  is the embedding vector of the choice-permuted questions that the model answered incorrectly, and  $\mathbf{z}_+$  is from the correctly answered questions. Also,  $n_-$  and  $n_+$  correspond to the number

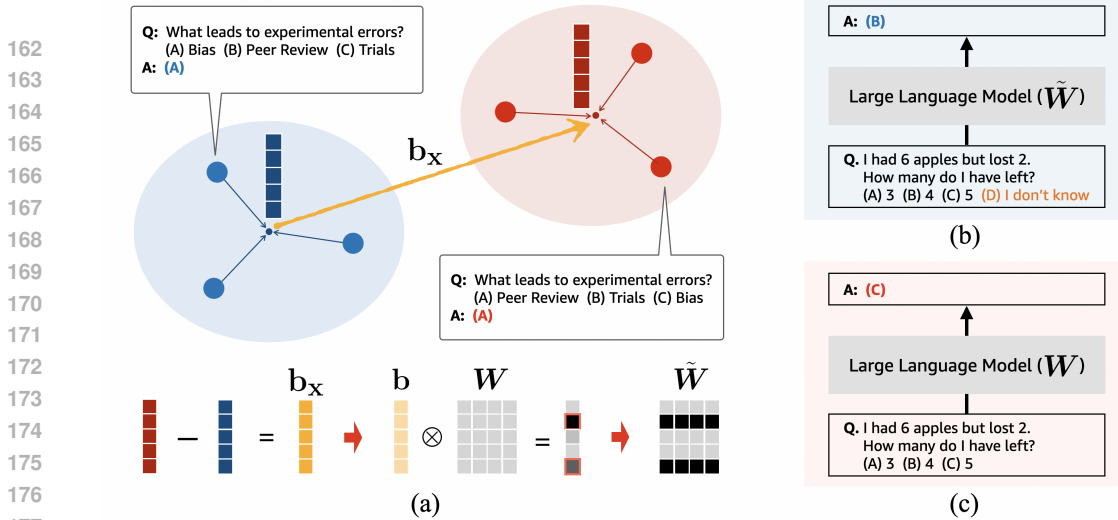


Figure 4: **Bias Node Pruning with Auxiliary Option Injection.** (a) The bias vector  $\mathbf{b}_x$  is computed for each sample using its choice-permuted embeddings (equation 1). The bias vectors are averaged across a small subset of training data to retrieve the average bias vector,  $\mathbf{b}$  (equation 2). Then,  $\mathbf{b}$  is used to select nodes to prune in  $\mathbf{W}$ , where  $\otimes$  refers to the operation in equation 4. (b) The pruned  $\tilde{\mathbf{W}}$  is used to retrieve answers for the test questions, along with our Auxiliary Option Injection technique that injects the “I don’t know” option in the inputs (§ 3.2). Our debiasing approaches may correct potentially erroneous responses retrieved with  $\mathbf{W}$  and without AOI, as in (c).

of incorrect and correct questions, respectively. To balance the number of correct and incorrect samples, we use the vector sets  $\{\mathbf{z}_-, \mathbf{z}_+\}$  only when  $1 \leq n_+/n_- \leq 2$ . Then, we average the bias vectors across the samples in data subset  $\mathcal{X}$  to define the average bias vector

$$\mathbf{b} = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}} \mathbf{b}_x, \quad (2)$$

where we use a subset size  $|\mathcal{X}|$  of 32 in this work. Refer to Figure 4(a) for visual aid.

We use the L2 norm of the average bias vector retrieved from different layers and tokens as a proxy for the magnitude of selection bias. Figure 3(b) shows the norm value from each location as a heatmap, where the x-axis lists the layer indices, and the y-axis shows the last 50 token embeddings of the inputs. Interestingly, the magnitude of the bias vector is prominent only in the final layer, motivating us to focus on the interaction of the average bias vector with the linear output head.

### 3 METHODS

Motivated by our observations that the selection bias is (1) prominently seen when the model is wrong, and (2) captured in the final decoder layers, we introduce two methods for debiasing the model predictions: **Bias Node Pruning** (BNP) and **Auxiliary Option Injection** (AOI). As the names suggest, BNP drops nodes in the final output layer that contribute to the selection bias, and AOI utilizes an auxiliary “I don’t know” option to eliminate bias induced by ignorance.

#### 3.1 BIAS NODE PRUNING

As shown in § 2.2, the average bias vector  $\mathbf{b} \in \mathbb{R}^d$  is most prominent in the final layer, and the selection bias materializes in the final output projection parameters,  $\mathbf{W} \in \mathbb{R}^{d \times |\mathcal{V}|}$ , where  $\mathcal{V}$  is the vocabulary set. To mitigate the selection bias problem induced by the linear layer, we prune the parameters in  $\mathbf{W}$  that contribute to the bias. In choosing which parameters to prune, we gain intuition by approximating a biased model,  $\mathcal{F}$ , as

$$\mathcal{F}(\mathbf{x}_A) \approx (\mathcal{D}(\mathbf{x}_A) + \mathbf{b}) \cdot \mathbf{W}, \quad (3)$$

where  $\mathcal{D}$  is a conceptual LLM decoder with zero selection bias, and  $\mathbf{b}$  is the average bias vector defined in equation 2. Then,  $\mathbf{b} \cdot \mathbf{W}$  is the factor that contributes to the selection bias, and removing the parameters in  $\mathbf{W}$  that has the most active interaction with  $\mathbf{b}$  will reduce selection bias. Accordingly, we choose the top- $k$  rows in  $\mathbf{W}$  with respect to

$$\mathcal{K} = \text{Top-}k \left( \sum_{i \in [1, d]} \mathbf{b}_i \times \mathbf{W}_{ij} \right), \quad (4)$$

where  $|\mathcal{V}|$  is the vocabulary size of the output. Then, we use the index in  $\mathcal{K}$  to zero out the corresponding rows (*i.e.*, nodes) in  $\mathbf{W}$ . Bias Node Pruning (BNP) is a one-time process with the average bias vector  $\mathbf{b}$  being pre-computed, and the pruned weight  $\tilde{\mathbf{W}}$  is applied to all test samples as  $f(\mathbf{x}_{\mathcal{A}}) \cdot \tilde{\mathbf{W}}$  where  $f$  is the LLM decoder. Refer to Appendix B for complexity analysis. Another design choice would deduct  $\mathbf{b}$  from the decoder output embedding; however, we observed more stable performance by pruning the parameters in  $\mathbf{W}$ .

### 3.2 AUXILIARY OPTION INJECTION

Because selection bias is more likely when a model is incorrect, we hypothesized that providing an ‘‘I don’t know’’ (IDK) option would reduce selection bias. The auxiliary option  $o_{\text{aux}}$  is applied as

$$\mathcal{A} := \mathcal{A} \cup \{o_{\text{aux}}\} \quad (5)$$

$$\hat{\mathbf{a}} = \arg \max_{a \in \mathcal{A} \setminus o_{\text{aux}}} P(\hat{\mathbf{y}} = a | \mathbf{x}_{\mathcal{A}}), \quad (6)$$

where  $\mathcal{A}$  is the set of answer choices, and  $\mathbf{x}_{\mathcal{A}}$  is the input question with choices  $\mathcal{A}$ . How we retrieve the probability for each choice  $a$  will be later discussed in the implementation details in § 5 and Appendix A.2. Further analyses on AOI will be provided in § 6.2.

## 4 EVALUATION

There is no consensus in the literature on how to measure selection bias. Here, we first review two selection bias metrics, Standard Deviation of Recalls (RStd) and Relative Standard Deviation (RSD), which evaluate the consistency of *performance* across choices. By scrutinizing their limitations, we propose **Choice Kullback-Leibler Divergence** (CKLD), which is a novel *distribution*-based bias metric.

**Definition 1.** (*Standard Deviation of Recalls*) is the standard deviation of the class-wise recall:

$$\text{RStd} = \sqrt{\frac{1}{k} \sum_{i=1}^k (r_i - \bar{r})^2}, \quad (7)$$

where  $k$  is the number of choices,  $r_i$  is the recall of the  $i$ -th class, and  $\bar{r}$  is the arithmetic mean of  $r_i$  values (Zheng et al., 2024).

**Definition 2.** (*Relative Standard Deviation*) is the class-wise accuracy standard deviation normalized by the overall accuracy:

$$\text{RSD} = \frac{\sqrt{\frac{1}{k} \sum_{i=1}^k (s_i - \bar{s})^2}}{\bar{s}}, \quad (8)$$

where  $k$  is the number of choices,  $s_i$  is the accuracy of the  $i$ -th class, and  $\bar{s}$  is the mean accuracy averaged across classes (Croce et al., 2021; Reif & Schwartz, 2024).

We empirically show how these performance-based metrics, RStd and RSD, behave across different data characteristics. We constructed synthetic 4-way MCQ datasets by varying the choice selection ratio under different ground-truth ratios. For instance, in the third column of Figure 5, labeled ‘‘A’ Label Ratio = 0.55’’, answer choice ‘A’ is the correct choice in 55% of the samples and the rest are labeled ‘B’, ‘C’, or ‘D’ 15% of the time, respectively. To simulate realistic predictions, we have the model render correct predictions half of the time, and predict with respect to the choice selection ratio (*i.e.*, ‘A’ selection rate) for the other half. For example, if ‘A’ Selection Rate is 0.4, each choice will be sampled with respect to  $P(A) = 0.4$  and  $P(B) = P(C) = P(D) = 0.2$  half of the time,

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

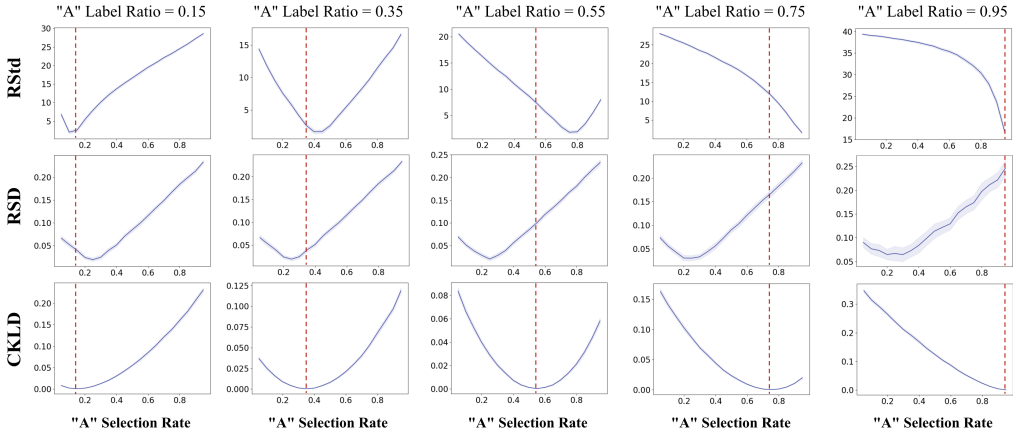


Figure 5: **Empirical analyses of selection bias metrics.** The metrics are tested on a 4-way classification task using synthetic data with varying levels of label ratios (outer x axis) and selection rates (inner x axis). We randomly generate 3000 samples and run 100 times to retrieve the mean and standard deviation of the metrics. The corresponding ‘A’ Ratios are denoted with dashed lines.

and will predict the correct answer for the other half. With this set up, the selection bias metrics should be lowest at the ‘A’ Label Ratio, shown with a vertical dashed line in Figure 5.

In contrast, the minimum points of RStd and RSD are not in the expected locations (Figure 5). Both metrics are insensitive to the ground-truth ratios. (RSD is lowest when the ‘A’ Selection Rate is  $\frac{1}{\# \text{ Choices}} = \frac{1}{4}$  regardless of the ‘A’ Label Ratio.) These results highlight the inability of RStd and RSD to measure selection bias in datasets with skewed distributions of the correct label. Therefore, we propose Choice Kullback-Leibler Divergence (CKLD), a distribution-based metric sensitive to data distribution and imbalance of choice labels.

**Definition 3.** (Choice Kullback-Leibler Divergence) is the KL divergence between the ratio of each predicted choice and the ratio of each ground truth choice label:

$$\text{CKLD} = \sum_{i=1}^k p_i \log \frac{p_i}{q_i}, \tag{9}$$

where  $k$  is the number of choices,  $p_i$  is the ratio of ground truth label choices, and  $q_i$  is the ratio of each predicted choice label.

CKLD is minimized when the predictions match the ground-truth ratio without bias towards certain choices (bottom row of Figure 5; proof in Appendix C and further discussion in Appendix C.1). However, CKLD does not account for the model performance in downstream tasks. Hence, it is important to refer to multiple metrics for a robust assessment. In this work, we leverage both RSD and our CKLD metrics to evaluate selection bias. We chose RSD because the groundtruth ratios of the benchmark datasets are close to uniform, and we can expect RSD to be minimized when the predictions are uniform.

## 5 EXPERIMENTS

In this section, we evaluate our Bias Node Pruning (BNP) and Auxiliary Option Injection (AOI) in various settings. We demonstrate the effect of our methods in § 5.1 and show that AOI can debias black-box models in § 5.2.

**Datasets and Models.** We evaluate our method on three multiple-choice question answering data test sets, ARC-Challenge (Clark et al., 2018), MMLU-Redux (Gema et al., 2024), and CommonsenseQA (Talmor et al., 2019). To retrieve the average bias vectors (equation 2), a separate set of out-of-bag samples is used. Further dataset details are provided in Appendix A.1. For the models, we mainly evaluate our approach on Llama-3-8B-Instruct (Meta, 2024), Mistral-7B-Instruct-v0.2 (Jiang et al., 2023), and Bloomz-7b1 (Muennighoff et al., 2023).

**Implementation Details.** As discussed in Section § 4, we employ RSD and CKLD to measure selection bias and assess the debiasing performance of our approach. We use Accuracy and the

Table 1: Bias Node Pruning (BNP) and Auxiliary Option Injection (AOI) are tested on three datasets with Llama-3, Bloomz, and Mistral. The best performances are in **bold**.

Method	ARC-Challenge				MMLU-Redux				CSQA			
	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓
Llama-3	52.3	54.1	0.562	0.494	41.8	46.7	1.021	0.589	65.4	66.2	0.261	0.095
Llama-3 + BNP	56.7	57.0	0.434	0.302	43.1	47.2	0.965	0.501	66.6	66.8	0.218	0.074
Llama-3 + AOI	60.7	61.0	0.364	0.231	47.3	49.9	0.807	0.321	67.4	67.8	0.211	0.065
Llama-3 + BNP + AOI	<b>65.3</b>	<b>65.1</b>	<b>0.262</b>	<b>0.124</b>	<b>48.3</b>	<b>50.5</b>	<b>0.531</b>	<b>0.288</b>	<b>68.1</b>	<b>68.2</b>	<b>0.174</b>	<b>0.049</b>
Bloomz	43.9	44.2	0.461	0.283	28.0	32.8	1.003	0.661	58.5	57.2	0.215	0.136
Bloomz + BNP	46.8	47.0	0.352	0.191	31.0	33.0	<b>0.537</b>	0.326	61.4	60.9	0.178	0.083
Bloomz + AOI	<b>48.9</b>	48.5	0.590	0.147	29.5	32.7	0.808	0.456	64.2	63.6	<b>0.134</b>	0.060
Bloomz + BNP + AOI	48.8	<b>48.9</b>	<b>0.208</b>	<b>0.088</b>	<b>32.0</b>	<b>33.3</b>	0.672	<b>0.205</b>	<b>64.9</b>	<b>64.9</b>	0.159	<b>0.052</b>
Mistral	67.4	67.6	0.156	0.040	46.4	47.6	0.366	0.186	63.6	63.9	0.184	0.042
Mistral + BNP	67.2	67.3	0.157	0.040	46.4	47.6	0.366	0.186	63.7	64.0	0.180	0.041
Mistral + AOI	<b>69.8</b>	<b>69.9</b>	<b>0.108</b>	<b>0.019</b>	<b>48.6</b>	<b>49.3</b>	<b>0.308</b>	<b>0.139</b>	<b>66.8</b>	<b>66.8</b>	0.101	<b>0.016</b>
Mistral + BNP + AOI	69.5	69.5	<b>0.108</b>	<b>0.019</b>	<b>48.6</b>	<b>49.3</b>	0.309	0.140	<b>66.8</b>	<b>66.8</b>	<b>0.099</b>	<b>0.016</b>

Table 2: **Comparison with Baselines.** Ours (BNP + AOI) is compared and applied to baseline methods. Best performances are in **bold**, and values denoted with \* are Ours with only BNP. Note that Bloomz + DoLa performed poorly and was meaningless to compare with baselines.

Method	ARC-Challenge				MMLU-Redux				CSQA			
	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓
Llama-3	52.3	54.1	0.562	0.494	41.8	46.7	1.021	0.589	65.4	66.2	0.261	0.095
Llama-3 + Ours	65.3	65.1	0.262	0.124	48.3	50.5	0.531	0.288	68.1	68.2	0.174	0.049
Llama-3 + CoT	66.2	66.3	0.178	0.050	50.2	51.0	0.641	0.124	65.3	65.7	0.161	0.025
Llama-3 + CoT + Ours	69.2	69.5	<b>0.156</b>	<b>0.024</b>	<b>50.4</b>	<b>51.1</b>	<b>0.281</b>	<b>0.095</b>	65.9	66.0	0.123	<b>0.012</b>
Llama-3 + ICL	62.2	61.7	0.292	0.169	42.6	46.4	0.735	0.486	69.0	69.0	<b>0.116</b>	0.026
Llama-3 + ICL + Ours	<b>70.0</b>	<b>70.0</b>	0.167	0.054	46.9	49.2	0.526	0.280	<b>69.5</b>	<b>69.3</b>	0.124	0.037
Llama-3 + DoLa	51.1	52.8	0.578	0.524	41.5	46.3	1.033	0.581	65.1	65.6	0.244	0.087
Llama-3 + DoLa + Ours	64.1	63.7	0.271	0.139	47.6	49.8	0.545	0.292	66.7	66.7	0.178	0.052
Bloomz	43.9	44.2	0.461	0.283	28.0	32.8	1.003	0.661	58.5	57.2	0.215	0.136
Bloomz + Ours	48.8	48.9	0.208	0.088	32.0	33.3	0.672	0.205	<b>64.9</b>	<b>64.9</b>	0.159	0.052
Bloomz + CoT	47.5	47.2	0.169	0.070	30.7	32.2	0.445	0.162	62.7	62.6	<b>0.093</b>	<b>0.020</b>
Bloomz + CoT + Ours	<b>50.2</b>	<b>50.1</b>	<b>0.058</b>	<b>0.013</b>	<b>34.3</b>	<b>34.7</b>	<b>0.215</b>	<b>0.019</b>	62.8*	62.8*	0.104*	<b>0.020*</b>
Bloomz + ICL	39.9	42.2	0.534	0.298	30.4	32.0	0.566	0.272	50.3	52.1	0.434	0.239
Bloomz + ICL + Ours	42.8*	45.2*	0.433*	0.249*	30.7*	31.1*	0.310*	0.135*	55.5	57.3	0.365	0.167
Mistral	67.4	67.6	0.156	0.040	46.4	47.6	0.366	0.186	63.6	63.9	0.184	0.042
Mistral + Ours	<b>69.5</b>	<b>69.5</b>	0.108	0.019	48.6	49.3	0.309	0.140	<b>66.8</b>	66.8	0.099	0.016
Mistral + CoT	66.6	66.5	0.510	0.021	50.3	50.5	0.551	0.063	63.2	63.4	0.476	0.025
Mistral + CoT + Ours	66.9	66.8	<b>0.071</b>	<b>0.014</b>	<b>50.6</b>	<b>50.7</b>	0.527	<b>0.032</b>	64.5	64.5	0.127	0.021
Mistral + ICL	65.7	66.0	0.183	0.054	43.1	44.5	0.410	0.253	61.7	61.7	0.167	0.046
Mistral + ICL + Ours	65.7	65.7	0.127	0.032	44.6	45.8	0.382	0.203	63.4	63.5	0.118	0.026
Mistral + DoLa	67.4	67.5	0.155	0.040	46.4	47.6	0.363	0.184	63.6	63.9	0.184	0.042
Mistral + DoLa + Ours	69.4	69.4	0.106	0.019	48.7	49.4	<b>0.305</b>	0.135	<b>66.8</b>	<b>66.9</b>	<b>0.098</b>	<b>0.015</b>

weighted F1 score for question answering performance evaluation. In predicting the answers from LLMs, we follow previous works (Zheng et al., 2024): we select the choice symbol (e.g., A, B, C, D) with the highest probability. For BNP, we prune 32 nodes for Llama-3 and Mistral, and 128 nodes for Bloomz. Because we modify the inference step, the entire process is not stochastic. More detailed explanation and further implementation details are provided in Appendix A.2.

## 5.1 MAIN EXPERIMENTS

**BNP + AOI consistently improves base model performance by reducing selection bias.** Table 1 shows the performance of our methods with three LLMs and MCQ datasets. For all models and data sets, BNP and/or AOI increased accuracy and F1 score and decreased RSD and CKLD. It is especially noteworthy that Llama-3’s accuracy on ARC-Challenge improves from 52.3% to 65.3% when both BNP and AOI are applied; an outstanding 24.9% increase.

**Our method can be applied together with other debiasing and decoding methods.** For further insight, we compare our methods with other debiasing and decoding approaches: Chain-of-Thought (CoT; Wei et al. (2022)), In-Context Learning (ICL; Brown et al. (2020)), and Decoding by Contrasting Layers (DoLa; Chuang et al. (2023)). For CoT, we follow the implementation of OpenAI Evals (OpenAI) by first prompting with “Let’s think step by step”, and then using the generated explanation to regenerate the final prediction. In the case of ICL, we take one question from the training set to retrieve  $N!$  choice-permuted questions, where  $N$  is the number of choices. Then, we randomly select three questions from the choice-permuted pool and create demonstrative examples from them, where the LLM agent always answers the choice-permuted questions correctly. Concrete prompt formats and details are provided in Appendix A.3. These baseline methods can be

Table 3: **Applying AOI to black-box settings.** For Llama-3, Bloomz, and Mistral, we assume that we do not have access to the parameters nor the probability outputs, identical to black-box models.

Method	ARC-Challenge				MMLU-Redux				CSQA			
	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓
Llama-3	65.7	65.8	0.086	0.007	51.9	52.2	0.184	0.034	69.9	69.8	0.051	0.003
Llama-3 + AOI	66.9	66.9	0.076	0.007	52.6	53.0	0.177	0.033	71.3	71.2	0.030	0.003
Bloomz	41.9	42.6	0.703	0.208	27.6	31.0	1.102	0.523	55.9	55.3	0.252	0.142
Bloomz + AOI	44.7	45.0	0.305	0.155	29.4	31.8	0.972	0.413	59.2	58.2	0.180	0.105
Mistral	55.2	55.2	0.140	0.036	47.4	47.6	0.216	0.069	54.6	54.8	0.155	0.031
Mistral + AOI	59.0	59.0	0.117	0.020	48.5	48.8	0.217	0.069	62.8	62.8	0.082	0.013
Claude-3-Haiku	65.3	65.0	0.095	0.024	52.1	52.0	0.057	0.008	36.4	37.3	0.587	0.331
Claude-3-Haiku + AOI	71.4	71.5	0.087	0.004	51.7	51.7	0.052	0.004	47.0	47.9	0.302	0.023
Claude-3-Sonnet	86.9	86.9	0.034	0.001	60.6	60.7	0.133	0.024	71.0	70.8	0.072	0.015
Claude-3-Sonnet + AOI	87.6	87.6	0.027	0.001	60.3	60.4	0.111	0.019	73.1	72.7	0.057	0.022

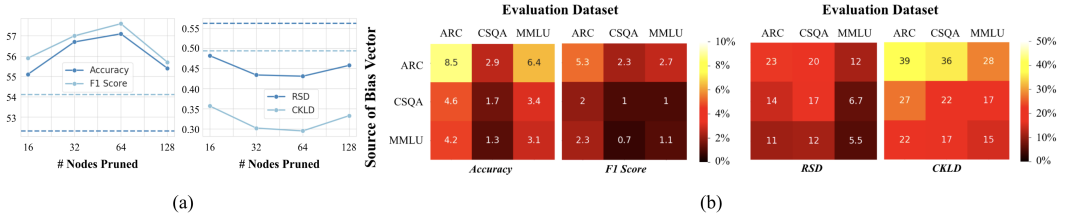


Figure 6: **BNP Analyses.** (a) BNP improves the base performances (dashed lines) regardless of the number of nodes pruned. The number of nodes to prune can be adjusted to achieve better performance. More figures are in Appendix D. (b) Each metric improvement (%) from its base Llama-3 performance when using the average bias vector from different sources is shown in heatmaps.

used along with our debiasing methods. Both question answering and debiasing improve when our methods are applied together (Table 2), even achieving the best performance in collaboration with appropriate baselines. Note that the values denoted with ‘\*’ are measured only when our BNP is applied because AOI did not fare well in those cases.

## 5.2 BLACK-BOX SETTINGS

Several state-of-the-art models are black-box and their parameters are not open to the public. In these cases, BNP is not feasible, leaving AOI as the only available technique for debiasing the model, using text outputs for prediction. For this reason, we devise a comparative experiment where only AOI is applied to the models. For white-box models Llama-3, Bloomz, and Mistral, we compute the Jaccard similarity between each choice option and the generated text to select the choice with the highest similarity score, instead of the probability-based answer selection method used in our main experiments. This approach simulates a black-box setting with the white-box models. Moreover, we extend our experiment to Claude-3 Haiku and Sonnet models (Anthropic, 2023), which are closed-source black-box models. In Table 3, AOI generally improves black-box model performance (accuracy and F1) and reduces selection bias (RDS and CKLD).

## 6 ANALYSES

In this section, we provide in-depth analyses on the mechanism and efficacy of our methods: Bias Node Pruning (§ 6.1) and Auxiliary Option Injection (§ 6.2). The qualitative findings from our experiments are discussed in § 6.3.

### 6.1 ANALYZING BIAS NODE PRUNING

**BNP is not sensitive to the number of nodes pruned.** Figure 6(a) reveals how the performance metrics change as the number of pruned nodes varies. Regardless of the number of nodes pruned from 8 to 128, our method improves the base performance (dashed lines in the figure) by great margins. While our method is robust to the amount of nodes pruned, searching for the adequate level of pruning may achieve better debiasing performance on the downstream task. Full list of the figure is provided in Appendix D.3.

**The average bias vector can be generalized across datasets.** The average bias vector represents the direction of selection bias in the embedding space. If the bias vector captures pure information



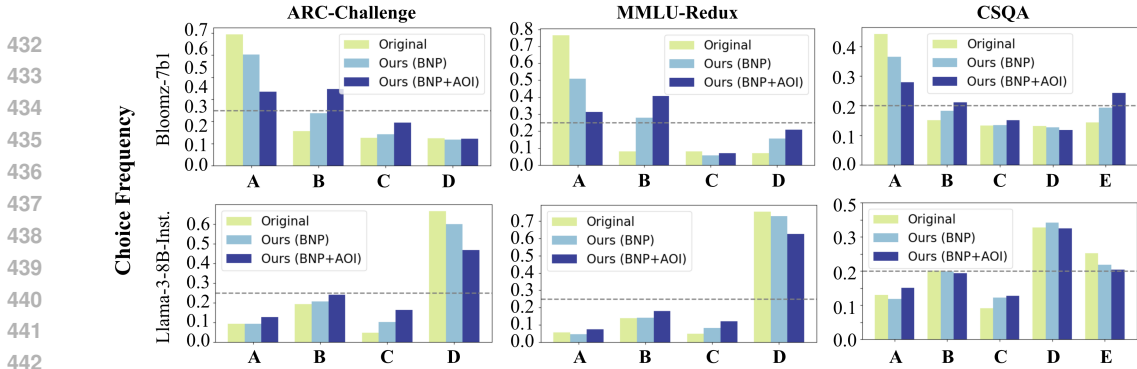


Figure 7: **Effect of our methods on choice distributions.** Our methods reduce the level of selection bias, and the choice distributions become flatter. Dashed lines are the uniform ratios (gold standard).

about selection bias, it should generalize across datasets. To test this hypothesis, we used the bias vector from one dataset on another. Figure 6(b) shows a heatmap of the improvement in each performance metric. Interestingly, there is no diagonal pattern, indicating bias vectors retrieved from one dataset can reduce selection bias in other datasets. For instance, the bias vector from the ARC-Challenge dataset improves the CKLD value of the CSQA dataset by 36%, which is even higher than the 22% improvement using the bias vector retrieved from its own CSQA dataset.

### 6.2 ANALYZING AUXILIARY OPTION INJECTION

#### Content of the auxiliary option matters.

Our experiments above used “I don’t know” as the auxiliary option, but other options are also possible. We conducted an experiment where we substituted it with “None of the above” and “I know the answer”. In Table 4, “None” refers to the former, and “Know” refers to the latter type of auxiliary option. For Llama-3 and Bloomz, the inclusion of an auxiliary option improves performance and reduces selection bias relative to the baseline (Table 4), but the “I don’t know” (Ours) performs better in most cases. With the Mistral model, however, the “I know the answer” option degrades model performance and increases selection bias. A full table with other datasets and more ablation experiments are in Appendix E

Table 4: AOI with different option contents on the MMLU-Redux dataset.

Method	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓
Llama-3	41.8	46.7	1.021	0.589
Llama-3 + "None"	42.4	42.7	0.833	0.487
Llama-3 + "know"	45.6	46.5	0.790	0.366
Llama-3 + Ours	<b>48.3</b>	<b>50.5</b>	<b>0.531</b>	<b>0.288</b>
Bloomz	28.0	32.8	1.003	0.661
Bloomz + "None"	26.5	25.9	0.730	0.518
Bloomz + "Know"	28.0	26.1	<b>0.618</b>	0.314
Bloomz + Ours	<b>32.0</b>	<b>33.3</b>	0.672	<b>0.205</b>
Mistral	46.4	47.6	0.366	0.186
Mistral + "None"	48.0	47.8	0.596	0.159
Mistral + "Know"	9.7	3.9	0.762	1.888
Mistral + Ours	<b>48.6</b>	<b>49.3</b>	<b>0.309</b>	<b>0.140</b>

### 6.3 QUALITATIVE EVALUATION

**Impact on choice distributions.** In Figure 7, we show how the distribution of the selected answer choices changes when we introduce BNP and AOI. In all three datasets, the distribution becomes more uniform when BNP and/or AOI are applied, indicating lower levels of selection bias. More qualitative examples are provided in Appendix F.

**Qualitative examples.** In addition to disclosing the distributional effect, we provide below the qualitative question-response examples of Llama-3 and Bloomz on the ARC-Challenge dataset. As in Figure 3(a), Llama-3 often showed a preference for choice ‘D’, regardless of the order of choices. Our method successfully corrects such errors. Bloomz, on the other hand, showed a preference for choice ‘A’. Again, our methods corrected the model’s response.

<p><b>Original Question:</b> Which of the following organs in fish has the same function as the human lung? (A) kidney (B) heart (C) skin (D) gill</p> <p>⇒ <b>Llama-3 Response:</b> (D)</p>	<p><b>Ground-truth:</b> (D)</p>
--	---------------------------------

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

**Permuted Question:** Which of the following organs in fish has the same function as the human lung? (A) kidney (B) heart (C) gill (D) skin

⇒ **Llama-3 Response:** (D) / **BNP+AOI Response:** (C)

**Ground-truth:** (C)

**Original Question:** Cells take in food for energy. The part of the cell that aids in the digestion of the food is the lysosome. What is the main role of lysosomes in the process of food digestion? (A) breaking down wastes (B) building proteins (C) controlling the activities of the cell (D) converting energy from one form into another

⇒ **Bloomz Response:** (A)

**Ground-truth:** (A)

**Permuted Question:** Cells take in food for energy. The part of the cell that aids in digestion of the food is the lysosome. What is the main role of lysosomes in the process of food digestion? (A) building proteins (B) breaking down wastes (C) controlling the activities of the cell (D) converting energy from one form into another

⇒ **Bloomz Response:** (A) / **BNP+AOI Response:** (B)

**Ground-truth:** (B)

## 7 RELATED WORKS

**Selection Bias.** The large language models’ tendency to favor choices in a certain order or with a specific symbol has been discussed in many previous works. Some of the works investigated the skewed pattern of responses for MCQs (Zheng et al., 2024; Wei et al., 2024; Pezeshkpour & Hruschka, 2024), emphasizing that selection bias is a critical problem. Many works have approached this problem by calibrating the output probabilities (Wang et al., 2023; Zheng et al., 2024; Reif & Schwartz, 2024; Wei et al., 2024; Pezeshkpour & Hruschka, 2024; Wang et al., 2024; Balepur et al., 2024; Li & Gao, 2024; Gupta et al., 2024), while others change the way queries are input (Li et al., 2023b; Robinson et al., 2023). Additional approaches include debiasing the LLM through distillation training (Liusie et al., 2024) and training the model to enforce its multiple choice symbol binding (MCSB) property (Xue et al., 2024). While parameter pruning methods are often used for efficient deep learning (Srinivas & Babu, 2015; Han et al., 2016; Zhu & Gupta, 2017; Molchanov et al., 2019; 2022) or to have the LLM unlearn certain factual knowledge (Liu et al., 2024; Pochinkov & Schoots, 2024), parameter pruning has rarely been discussed for debiasing. Thus, our Bias Node Pruning is a novel approach in the context of the selection bias.

**Auxiliary Options.** Inclusion of the “I don’t know” option can improve the quality of data collected in surveys (Schuman & Presser, 1996) but does not meaningfully impact the labels assigned by annotators (Beck et al., 2022). Recent research has drawn attention to the similarities between surveys, labeling tasks, and model responses to MCQs (Tjumatja et al., 2023; Eckman et al., 2024; Chen et al., 2024) Further research into LLMs’ response behavior would benefit from incorporating insights from the survey science domain: see the discussion in (Eckman et al., 2024).

## 8 CONCLUSION

When LLMs answer MCQs, selection bias is a critical problem. Previous research has predominantly focused on modifying the LLM’s input and/or output. In contrast, we uncover the *internal* source of the bias by scrutinizing the embedding-level discrepancies introduced by this bias. Building on these insights, we propose Bias Node Pruning (BNP) and Auxiliary Option Injection (AOI). Additionally, we address the limitations of existing performance-based evaluation metrics by introducing a new distribution-based metric, Choice Kullback-Leibler Divergence (CKLD), which addresses the insensitivity of prior metrics to imbalance of choice labels. Our approach improved MCQ answering performance by reducing the level of selection bias across widely used MCQ datasets using both open-source (white box) and closed-source (black-box) models. BNP and AOI work alongside other debiasing/decoding methods to improve the base performance of Llama-3 by up to 33.8% on the ARC-Challenge dataset. We also conducted in-depth analyses to better understand the effect of each component, along with case studies to provide qualitative insight. Overall, our method provides a novel intuition in scrutinizing the internal source of selection bias, and also provides a new approach in debiasing LLMs.

## REFERENCES

- 540  
541  
542 Anthropic. Introducing claude, 2023. URL [https://www.anthropic.com/news/  
543 introducing-claude](https://www.anthropic.com/news/introducing-claude).
- 544 Nishant Balepur, Abhilasha Ravichander, and Rachel Rudinger. Artifacts or abduction: How do  
545 llms answer multiple-choice questions without the question? *arXiv preprint arXiv:2402.12483*,  
546 2024.
- 547  
548 Jacob Beck, Stephanie Eckman, Rob Chew, and Frauke Kreuter. Improving labeling through social  
549 science insights: results and research agenda. In *International Conference on Human-Computer  
550 Interaction*, pp. 245–261. Springer, 2022.
- 551  
552 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
553 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
554 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 555  
556 Nuo Chen, Jiqun Liu, Xiaoyu Dong, Qijiong Liu, Tetsuya Sakai, and Xiao-Ming Wu. Ai can be  
557 cognitively biased: An exploratory study on threshold priming in llm-based batch relevance as-  
558 sessment. *arXiv preprint arXiv:2409.16022*, 2024.
- 559  
560 Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R Glass, and Pengcheng He. Dola:  
561 Decoding by contrasting layers improves factuality in large language models. In *The Twelfth  
562 International Conference on Learning Representations*, 2023.
- 563  
564 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
565 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
566 *arXiv preprint arXiv:1803.05457*, 2018.
- 567  
568 Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo Debenedetti, Nicolas Flam-  
569 marion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversar-  
570 ial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems  
571 Datasets and Benchmarks Track (Round 2)*, 2021.
- 572  
573 Stephanie Eckman, Barbara Plank, and Frauke Kreuter. Position: Insights from survey methodology  
574 can improve training data. In *Forty-first International Conference on Machine Learning*, 2024.
- 575  
576 Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria  
577 Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani,  
578 et al. Are we done with mmlu? *arXiv preprint arXiv:2406.04127*, 2024.
- 579  
580 Vipul Gupta, David Pantoja, Candace Ross, Adina Williams, and Megan Ung. Changing answer  
581 order can decrease mmlu accuracy. *arXiv preprint arXiv:2406.19470*, 2024.
- 582  
583 Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks  
584 with pruning, trained quantization and huffman coding. In *International Conference on Learning  
585 Representations*, 2016.
- 586  
587 Xingwei He, Zhenghao Lin, Yeyun Gong, A-Long Jin, Hang Zhang, Chen Lin, Jian Jiao, Siu Ming  
588 Yiu, Nan Duan, and Weizhu Chen. AnnoLLM: Making large language models to be better crowd-  
589 sourced annotators. In *Proceedings of the 2024 Conference of the North American Chapter of the  
590 Association for Computational Linguistics: Human Language Technologies (Volume 6: Industry  
591 Track)*, pp. 165–190, 2024.
- 592  
593 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*, 2021.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.

- 594 Minzhi Li, Taiwei Shi, Caleb Ziems, Min-Yen Kan, Nancy Chen, Zhengyuan Liu, and Diyi Yang.  
595 Coannotating: Uncertainty-guided work allocation between human and large language models  
596 for data annotation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural*  
597 *Language Processing*, pp. 1487–1505, 2023a.
- 598  
599 Ruizhe Li and Yanjun Gao. Anchored answers: Unravelling positional bias in gpt-2’s multiple-  
600 choice questions. *arXiv preprint arXiv:2405.03205*, 2024.
- 601  
602 Zongjie Li, Chaozheng Wang, Pingchuan Ma, Daoyuan Wu, Shuai Wang, Cuiyun Gao, and Yang  
603 Liu. Split and merge: Aligning position biases in large language model based evaluators. *arXiv*  
604 *preprint arXiv:2310.01432*, 2023b.
- 605  
606 Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, PRANAY SHARMA, Sijia  
607 Liu, et al. Model sparsity can simplify machine unlearning. *Advances in Neural Information*  
*Processing Systems*, 36, 2024.
- 608  
609 Adian Liusie, Yassir Fathullah, and Mark JF Gales. Teacher-student training for debiasing: General  
610 permutation debiasing for large language models. *arXiv preprint arXiv:2403.13590*, 2024.
- 611  
612 Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024. URL  
613 <https://ai.meta.com/blog/meta-llama-3/>.
- 614  
615 Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. Importance estimation  
616 for neural network pruning. In *Proceedings of the IEEE/CVF conference on computer vision and*  
*pattern recognition*, pp. 11264–11272, 2019.
- 617  
618 Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional  
619 neural networks for resource efficient inference. In *International Conference on Learning Repre-*  
*sentations*, 2022.
- 620  
621 Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven  
622 Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, et al. Crosslingual  
623 generalization through multitask finetuning. In *Proceedings of the 61st Annual Meeting of the*  
624 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15991–16111, 2023.
- 625  
626 OpenAI. Openai evals. <https://github.com/openai/evals>.
- 627  
628 Pouya Pezeshkpour and Estevam Hruschka. Large language models sensitivity to the order of op-  
629 tions in multiple-choice questions. In *Findings of the Association for Computational Linguistics:*  
*NAACL 2024*, pp. 2006–2017, 2024.
- 630  
631 Nicholas Pochinkov and Nandi Schoots. Dissecting language models: Machine unlearning via  
632 selective pruning. *arXiv preprint arXiv:2403.01267*, 2024.
- 633  
634 Yuval Reif and Roy Schwartz. Beyond performance: Quantifying and mitigating label bias in llms.  
635 In *Proceedings of the 2024 Conference of the North American Chapter of the Association for*  
636 *Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6784–  
6798, 2024.
- 637  
638 Joshua Robinson, Christopher Michael Rytting, and David Wingate. Leveraging large language  
639 models for multiple choice question answering. In *ICLR*, 2023.
- 640  
641 Howard Schuman and Stanley Presser. Questions and answers in attitude surveys: experiments on  
642 question forms, wording, and context, 1996.
- 643  
644 Robyn Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: An open multilingual graph of  
645 general knowledge. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31,  
2017.
- 646  
647 Suraj Srinivas and R Venkatesh Babu. Data-free parameter pruning for deep neural networks. In  
*Proceedings of the British Machine Vision Conference 2015*. British Machine Vision Association,  
2015.

- 648 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A ques-  
649 tion answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Con-*  
650 *ference of the North American Chapter of the Association for Computational Linguistics: Human*  
651 *Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, 2019.
- 652  
653 Lindia Tjuaaja, Valerie Chen, Sherry Tongshuang Wu, Ameet Talwalkar, and Graham Neubig.  
654 Do llms exhibit human-like response biases? a case study in survey design. *arXiv preprint*  
655 *arXiv:2311.04076*, 2023.
- 656 Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu,  
657 Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators. *arXiv preprint*  
658 *arXiv:2305.17926*, 2023.
- 659  
660 Xinpeng Wang, Bolei Ma, Chengzhi Hu, Leon Weber-Genzel, Paul Röttger, Frauke Kreuter, Dirk  
661 Hovy, and Barbara Plank. ” my answer is c”: First-token probabilities do not match text answers  
662 in instruction-tuned language models. *arXiv preprint arXiv:2402.14499*, 2024.
- 663 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
664 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *Advances*  
665 *in neural information processing systems*, volume 35, pp. 24824–24837, 2022.
- 666  
667 Sheng-Lun Wei, Cheng-Kuang Wu, Hen-Hsen Huang, and Hsin-Hsi Chen. Unveiling selection  
668 biases: Exploring order and token sensitivity in large language models. In *ACL*, 2024.
- 669 Mengge Xue, Zhenyu Hu, Meng Zhao, Liquan Liu, Kuo Liao, Shuang Li, Honglin Han, and Cheng-  
670 guo Yin. Strengthened symbol binding makes large language models reliable multiple-choice  
671 selectors. In *ACL*, 2024.
- 672  
673 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
674 chine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association*  
675 *for Computational Linguistics*, pp. 4791–4800, 2019.
- 676 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models  
677 are not robust multiple choice selectors. In *The Twelfth International Conference on Learning*  
678 *Representations*, 2024.
- 679 Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for  
680 model compression. *arXiv preprint arXiv:1710.01878*, 2017.
- 681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

# Appendix

## Table of Contents

<b>A Further Experimental Details</b>	<b>14</b>
A.1 Datasets . . . . .	14
A.2 Implementation Details . . . . .	15
A.3 Baselines . . . . .	15
A.4 Metrics . . . . .	16
<b>B Complexity of Bias Node Pruning</b>	<b>18</b>
<b>C Proof of CKLD’s label ratio sensitivity</b>	<b>18</b>
C.1 Why does an LLM need to match the ground truth ratio? . . . . .	19
<b>D More Experiments and Analyses</b>	<b>19</b>
D.1 Significance Test . . . . .	19
D.2 Further experiments on HellaSwag dataset . . . . .	19
D.3 Extended List of Figures . . . . .	19
<b>E Different AOI Setup</b>	<b>20</b>
<b>F Qualitative Examples</b>	<b>21</b>
<b>G Limitations and Broader Impact</b>	<b>26</b>

## A FURTHER EXPERIMENTAL DETAILS

### A.1 DATASETS

We experiment on three datasets: ARC-Challenge (Clark et al., 2018), MMLU-Redux (Gema et al., 2024), and CommonsenseQA Talmor et al. (2019). We also provide the ground-truth choice ratios in the test dataset in Table 5.

**ARC-Challenge** is a dataset from the AI2 Reasoning Challenge, which contains grade-school level multiple-choice science questions. Among the ‘Challenge’ and the ‘Easy’ sets, we use the former set with 1.17K test and 1.12K training questions. The training questions are used to extract the average bias vectors.

**MMLU-Redux** is a dataset derived from the original Massive Multitask Language Understanding (MMLU) (Hendrycks et al., 2021) dataset, which comprises multiple-choice questions from 57 different branches of knowledge. Gema et al. (2024) discovered that this original version contains numerous errors, and curated the dataset to have 3,000 manually re-annotated questions across 30 subjects in the original MMLU dataset. In the case of MMLU-Redux, there is no training set available. So we utilize the validation set from the original MMLU dataset to pre-compute the average bias vectors.

**CommonsenseQA** is a dataset of multiple-choice questions that require commonsense knowledge to respond. The dataset questions are extracted using the knowledge graph, ConceptNet (Speer et al., 2017), which consists of 9.74K training and 1.22K validation questions. We use the training set to retrieve the average bias vectors and evaluate on the validation set.

Table 5: **Ground-truth Label ratios of each dataset.**

Datasets	A ratio	B ratio	C ratio	D ratio	E ratio
ARC-Challenge	22.4%	25.7%	25.9%	24.1%	-
MMLU-Redux	22.3%	24.6%	25.4%	27.7%	-
CSQA	19.6%	20.9%	19.7%	20.6%	19.2%

## A.2 IMPLEMENTATION DETAILS

Here, we detail how we retrieve model predictions and list hyperparameters used for each model-dataset experiment.

**How are predictions retrieved?** As discussed in the main paper, we use the token output probability distribution to select a token ID for prediction. For instance, if  $\mathbf{z} \in \mathbb{R}^{|\mathcal{V}|}$  is the output logit vector of the first output token, we use  $\mathbf{z}[\text{'A'}] + \mathbf{z}[\text{'_A'}]$  to retrieve the logit value for choice 'A', and do the same for other choices as well. Note that '\_A' is a token that represents "A" with a space in front of it, whereas 'A' is a one-character token. Since these two represent the same choice, we aggregate their logits,  $\mathbf{z}$ , for accurate evaluation. Then, we take the softmax over all the choice logits to retrieve the final probability distribution over the choices.

**System prompt.** We use the same system prompt across all experiments: "You are an AI assistant that answers multiple choice questions. Please respond with capitalized alphabet(s) that correspond to the correct answer". For Chain-of-Thought reasoning baseline experiments, we use a slightly different version of "You are an AI assistant that answers multiple choice questions. Please think step by step and respond with capitalized alphabet(s) that correspond to the correct answer" to encourage the model to output a step-by-step reasoning process.

**Hyperparameters.** The number of nodes pruned is the main hyperparameter of our experiments. As disclosed in the main paper, we pruned 32 nodes in all experiments with Llama-3 and Mistral, and pruned 128 nodes in experiments with Bloomz. We did a simple hyperparameter search among {16, 32, 64, 128} nodes. Results can be found in Figure 6(a) and Figure 9. Another noteworthy hyperparameter is the choice delimiter, which refers to the type of token used to separate choices. In our preliminary experiments, we found that different choice delimiters such as space (' '), line break tokens ('\n'), multiple lines ('\n\n'), or special tags ('<c>') have varying impact on performance. As there were no consistent results, however, we chose to use the basic space delimiter in all our experiments, e.g. 'What is 1 + 1? (A) 2 (B) 3 (C) 4'. Although we do not discuss this in depth as it is beyond the scope of our work, we believe that analyzing the effect of different choice delimiters in multiple choice question answering would introduce an interesting viewpoint.

## A.3 BASELINES

In this section, we provide further details on how the debiasing baselines in Table 2 are designed.

**Chain-of-Thought (CoT)** first generates the model response that includes explanations by prompting with "Let's think step by step" as follows.

**System Prompt:** You are an AI assistant that answers multiple choice questions. Please think step by step and respond with capitalized alphabet(s) that correspond to the correct answer.

**User:** { *question* }.

**Assistant:** Let's think step by step.

Using the explanation that is generated with the prompt, we query the LLM once more with

**System Prompt:** You are an AI assistant that answers multiple choice questions. Please think step by step and respond with capitalized alphabet(s) that correspond to the correct answer.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

**User:** { *question* }.

**Assistant:** Let's think step by step. { *explanation* }. So the correct answer is

and identically use the first token output probability distribution to retrieve the predictions. Note that the actual prompt format depends on the model and the template above is a generic form.

**In-Context Learning (ICL)** takes one question out-of-bag sample and retrieve  $N!$  choice-permuted questions, where  $N$  is the number of choices. Then, three of the choice-permuted questions among the  $N!$  pool are randomly chosen to be used for the ICL demonstrative examples. Concretely, we design the prompt as follows.

**System Prompt:** You are an AI assistant that answers multiple choice questions. Please respond with capitalized alphabet(s) that correspond to the correct answer.

# Example 1

**User:** What leads to experimental errors? (A) Bias (B) Peer Review (C) Repeated Trials

**Assistant :** (A)

# Example 2

**User:** What leads to experimental errors? (A) Repeated Trials (B) Peer Review (C) Bias

**Assistant :** (C)

# Example 3

**User:** What leads to experimental errors? (A) Peer Review (B) Bias (C) Repeated Trials

**Assistant :** (B)

**User:** { *question* }.

**Assistant:**

Again, the prompt template is generic, and the actual input format depends on the model type.

**Decoding by Contrasting Layers (DoLa)** is a language model decoding method proposed by Chuang et al. (2023). Following their implementation, we measure the Jensen-Shannon Divergence between the final (or mature) output probability distribution and intermediate (or premature) outputs to select the layer with the highest divergence. Then, we use the selected layer output to divide the final output. Since this is similar to calibration, we expected DoLa to have debiasing effects. However, the results in Table 2 show that DoLa alone does not reduce the level of selection bias.

## A.4 METRICS

In this section, we provide a full list of selection bias metrics, including RStd, RSD, our CKLD, and other existing metrics that were not discussed in the main paper. We taxonomize the metrics into three groups: brute-force evaluation, performance-based evaluation, and distribution-based evaluation.

### A.4.1 BRUTE-FORCE EVALUATION

Brute-force evaluation metrics utilize all possible choice permutations to retrieve the metric value. Since we need to infer the output for each of the choice-permuted questions, the computation increases by a factor of  $N!$ , where  $N$  is the number of choices in the question. Here, we list two brute-force evaluation metrics, Proportion of Plurality Agreement (PPA) and Permutation Sensitivity (PS), and one semi-brute-force metric that additionally computes only the reverse-order permutation, Fluctuation Rate (FR).



**Definition 1.** (*Proportion of Plurality Agreement*) is the proportion of the plurality choice among all possible choice orderings of a multiple-choice question:

$$\text{PPA} = \frac{1}{|\mathcal{X}|} \sum_{\mathcal{X}} \frac{\max_n \left( \sum_{j=1}^N y_j = o_n \right)}{N!}, \quad (10)$$

where  $\mathcal{X}$  is the set of test samples,  $N$  is the number of choices in each question,  $n$  is the index of the choices,  $y_j$  is the choice content of the  $j$ -th choice-permuted sample prediction, and  $o_n$  is the  $n$ -th choice content. (Robinson et al., 2023)

**Definition 2.** (*Permutation Sensitivity*) is the expected divergence in output probability distributions of the choice-permuted questions:

$$\text{PS} = \mathbb{E}_{\sigma_i, j} [d(P(\cdot | q, \mathcal{A}_{\sigma_i}); P(\cdot | q, \mathcal{A}_{\sigma_j}))], \quad (11)$$

where  $\sigma_i$  is an arbitrary permutation of choices,  $\mathcal{A}_{\sigma_i}$  is the answer choice with the choice permutation,  $q$  is the input question,  $d(\cdot | \cdot)$  is the divergence function (e.g., KL-divergence), and  $P(\cdot | \cdot)$  is the output probability distribution function. (Liusie et al., 2024)

**Definition 3.** (*Fluctuation Rate*) is the rate of inconsistent model responses to the original input question and the question with choices presented in reversed order:

$$\text{FR} = \frac{1}{M} \sum_{i=1}^M \mathbf{1}(y_i \neq \overleftarrow{y}_i), \quad (12)$$

where  $M$  is the number of test questions,  $\mathbf{1}$  is the indicator function,  $\overrightarrow{y}$  is the model prediction to the original question, and  $\overleftarrow{y}$  is the prediction to the question with reversed choice order. (Wei et al., 2024)

#### A.4.2 PERFORMANCE-BASED EVALUATION

Performance-based evaluation tries to capture the consistency of model performance when measuring selection bias. The two metrics discussed in the paper, RStd and RSD, fall under this category.

**Definition 4.** (*Standard Deviation of Recalls*) is the standard deviation of the class-wise recall:

$$\text{RStd} = \sqrt{\frac{1}{k} \sum_{i=1}^k (r_i - \bar{r})^2}, \quad (13)$$

where  $k$  is the number of choices,  $r_i$  is the recall of the  $i$ -th class, and  $\bar{r}$  is the arithmetic mean of  $r_i$  values. (Zheng et al., 2024)

**Definition 5.** (*Relative Standard Deviation*) is the class-wise accuracy standard deviation normalized by the overall accuracy:

$$\text{RSD} = \frac{\sqrt{\frac{1}{k} \sum_{i=1}^k (s_i - \bar{s})^2}}{\bar{s}}, \quad (14)$$

where  $k$  is the number of choices,  $s_i$  is the accuracy of the  $i$ -th class, and  $\bar{s}$  is the mean accuracy averaged across classes. (Croce et al., 2021; Reif & Schwartz, 2024)

#### A.4.3 DISTRIBUTION-BASED EVALUATION

Existing performance-based evaluation metrics are insensitive to imbalance of choice labels, and manually adjusting the label distribution does not guarantee fair evaluation and may severely influence performance. Thus, we propose a new distribution-based evaluation metric, Choice Kullback-Leibler Divergence (CKLD), to complement evaluation of the selection bias.

**Definition 6.** (*Choice Kullback-Leibler Divergence*) is the KL divergence between the ratio of each predicted choice and the ratio of each ground truth choice label:

$$\text{CKLD} = \sum_{i=1}^k p_i \log \frac{p_i}{q_i}, \quad (15)$$

where  $k$  is the number of choices,  $p_i$  is the ratio of ground truth label choices, and  $q_i$  is the ratio of each predicted choice label.

## B COMPLEXITY OF BIAS NODE PRUNING

Bias Node Pruning is a two-step process that includes the (1) average bias vector computation, and (2) node pruning. The first phase utilizes  $M$  out-of-bag samples with  $N$  choices. This step requires computing the outputs of  $N!$  choice-permuted questions, translating to a complexity of  $O(N! \cdot M)$ . Once we retrieve the average bias vector, we use it to compute the top- $k$  nodes that activate selection bias (equation 4). This is also a one-time process whose node-pruned parameters are applied throughout all test-time inference tasks. The complexity of inference itself is identical to the original model without Bias Node Pruning, which is proportional to the number of test samples evaluated.

## C PROOF OF CKLD’S LABEL RATIO SENSITIVITY

We want to prove that CKLD is minimized when the prediction has no bias towards a certain choice, and matches the ratio of ground-truth labels. From the CKLD definition (equation 15) of

$$\text{CKLD} = \sum_{i=1}^k p_i \log \frac{p_i}{q_i}, \quad (16)$$

let  $q_i = p_i r_i$ , where  $r_i$  is the selection bias multiplier applied to the ground-truth choice ratio for each  $i = 1, \dots, k$ . As we want to find out when CKLD is minimized, we formulate the objective as follows:

$$\begin{aligned} & \text{minimize } \sum_{i=1}^k p_i \log \frac{p_i}{q_i} \\ & \text{s.t. } q_i = p_i r_i \text{ and } \sum_{i=1}^k p_i r_i = 1. \end{aligned} \quad (17)$$

By rewriting this as a Lagrangian function  $\mathcal{L}$ ,

$$\begin{aligned} \mathcal{L}(r_1, \dots, r_k, \lambda) &= \sum_{i=1}^k p_i \log \frac{p_i}{p_i r_i} + \lambda \left( \sum_{i=1}^k p_i r_i - 1 \right) \\ &= - \sum_{i=1}^k p_i \log r_i + \lambda \left( \sum_{i=1}^k p_i r_i - 1 \right), \end{aligned} \quad (18)$$

where  $\lambda$  is the Lagrangian multiplier, we take the partial derivative of each variable as:

$$\frac{\partial \mathcal{L}}{\partial r_i} = -\frac{p_i}{r_i} + \lambda p_i = 0 \quad (19)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{i=1}^k p_i r_i - 1 = 0. \quad (20)$$

Then, from equation 19,

$$r_i = \frac{1}{\lambda}, \quad (21)$$

and by substituting this to equation 20, we get

$$\begin{aligned} 0 &= \sum_{i=1}^k \frac{p_i}{\lambda} - 1 \\ &= \frac{1}{\lambda} - 1. \end{aligned} \quad (22)$$

Therefore, the objective is minimized when  $\lambda = 1$ , which translates to  $r_i = 1$  ( $\because$  equation 21). This is equivalent to saying that CKLD is minimized when  $q_i = p_i r_i = p_i$ , *i.e.*, when the prediction ratio matches the actual label ratio and there is no selection bias towards a certain choice.  $\square$

Table 6: Further experiments are done on the HellaSwag dataset.

	Acc	F1	RSD	CKLD
<i>ARC-Challenge</i>				
Llama-3	53.2 (1.3)	55.4 (1.3)	0.640 (0.142)	0.485 (0.049)
Llama-3 + BNP	57.4 (1.0)	58.0 (1.1)	0.533 (0.145)	0.304 (0.029)
Llama-3 + AOI	62.7 (1.0)	63.0 (1.1)	0.417 (0.133)	0.201 (0.023)
Llama-3 + BNP + AOI	66.8 (1.0)	66.6 (0.9)	0.340 (0.140)	0.121 (0.010)
<i>MMLU-Redux</i>				
Llama-3	39.8 (1.6)	44.4 (1.8)	0.982 (0.097)	0.673 (0.063)
Llama-3 + BNP	40.8 (1.7)	44.8 (1.8)	0.936 (0.100)	0.595 (0.065)
Llama-3 + AOI	44.5 (1.8)	47.0 (2.0)	0.657 (0.097)	0.384 (0.042)
Llama-3 + BNP + AOI	45.4 (1.6)	47.5 (1.8)	0.564 (0.018)	0.346 (0.041)
<i>CommonsenseQA</i>				
Llama-3	63.3 (1.1)	64.2 (0.9)	0.282 (0.026)	0.106 (0.018)
Llama-3 + BNP	64.9 (1.1)	65.2 (1.1)	0.222 (0.012)	0.073 (0.007)
Llama-3 + AOI	65.9 (0.9)	66.3 (0.8)	0.220 (0.020)	0.069 (0.010)
Llama-3 + BNP + AOI	67.2 (0.6)	67.2 (0.6)	0.175 (0.011)	0.052 (0.004)

### C.1 WHY DOES AN LLM NEED TO MATCH THE GROUND TRUTH RATIO?

Consider a scenario in which an LLM exhibits a bias toward selecting option ‘A’. In cases where the LLM is uncertain about the correct answer and resorts to random selection, it is more likely to choose ‘A’, resulting in a skewed overall choice distribution that diverges from the ground truth distribution. In contrast, an unbiased LLM would select options uniformly under uncertainty, producing a choice distribution that more closely aligns with the original ground truth distribution. Therefore, the extent to which an LLM’s predictions match the ground truth distribution can serve as a proxy for measuring Selection Bias.

## D MORE EXPERIMENTS AND ANALYSES

Here, we provide further experiments and analysis results that were not included in the main manuscript. In § D.2, we demonstrate an extended experiment result on another dataset. In § D.3, an extended list of figures of Figure 6 (a) is provided.

### D.1 SIGNIFICANCE TEST

In Table 6, we present the results of a significance test conducted on Llama-3 by performing 8 experiments, each with randomly permuted choices. The mean values for each dataset are reported, with standard deviations shown in parentheses. All values are statistically significant compared to the Llama-3 baseline, with t-test p-values below 0.001.

### D.2 FURTHER EXPERIMENTS ON HELLA SWAG DATASET

Beyond the three datasets tested in our main paper in Table 1, we disclose results on another widely used benchmark dataset, HellaSwag (Zellers et al., 2019). HellaSwag is a commonsense natural language inference (NLI) dataset that contains 4-way MCQ samples that asks the model to select the option that best ends the given sentence. The experimental results are in Table 7. Bloomz is not included in the table because the model failed to reasonably respond to most of the questions.

### D.3 EXTENDED LIST OF FIGURES

Here, we provide a comprehensive table of figures on the sensitivity test on the number of nodes pruned (§ 6.1, Figure 6(a)). In Figure 9, the effect of the number of pruned nodes is shown across the three models and datasets, as its value is varied from 16 to 128. We also provide the heatmap of the average bias vector magnitude in Figure 8. Similar to what has been shown in Figure 3 (b), selection bias seems prominent in the latter part of the decoder layers.

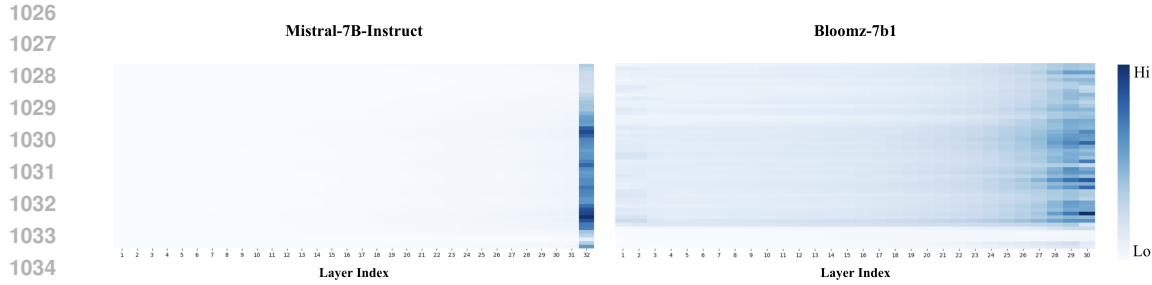


Figure 8: **More figures on different models other than Llama-3.** Left is the bias vector magnitude heatmap from Mistral-7B-Instruct, and right is from Bloomz-7b1.

Table 7: **Further experiments are done on the HellaSwag dataset.**

Method	HellaSwag			
	Acc. $\uparrow$	F1 $\uparrow$	RSD $\downarrow$	CKLD $\downarrow$
Llama-3	35.9	42.3	0.988	1.416
Llama-3 + BNP	38.6	43.6	0.861	0.998
Llama-3 + AOI	47.6	51.2	0.599	0.611
Llama-3 + BNP + AOI	<b>50.8</b>	<b>52.9</b>	<b>0.487</b>	<b>0.363</b>
Mistral	46.7	48.7	0.558	0.341
Mistral + BNP	46.5	48.6	0.563	0.345
Mistral + AOI	<b>51.7</b>	<b>53.0</b>	<b>0.414</b>	<b>0.206</b>
Mistral + BNP + AOI	51.6	52.9	0.415	0.207

## E DIFFERENT AOI SETUP

In this section, in addition to all three dataset ablation studies on the content of auxiliary options in § 6.2, we provide further ablation study results on the number and location of the auxiliary options.

**More auxiliary options have mixed effects on performance.** We find that controlling the number of auxiliary options has a notable impact on performance. That is, we tried adding multiple auxiliary options, all with the same “I don’t know” content. In most cases in Table 8, adding more auxiliary options did not help improve performance (see  $n$ -Choices AOI). Interestingly, however, both the question-answering and debiasing performance of Llama-3 significantly improved when using more options. This seems to be a peculiar property of Llama-3 that we can enhance its performance by simply adding multiple auxiliary options.

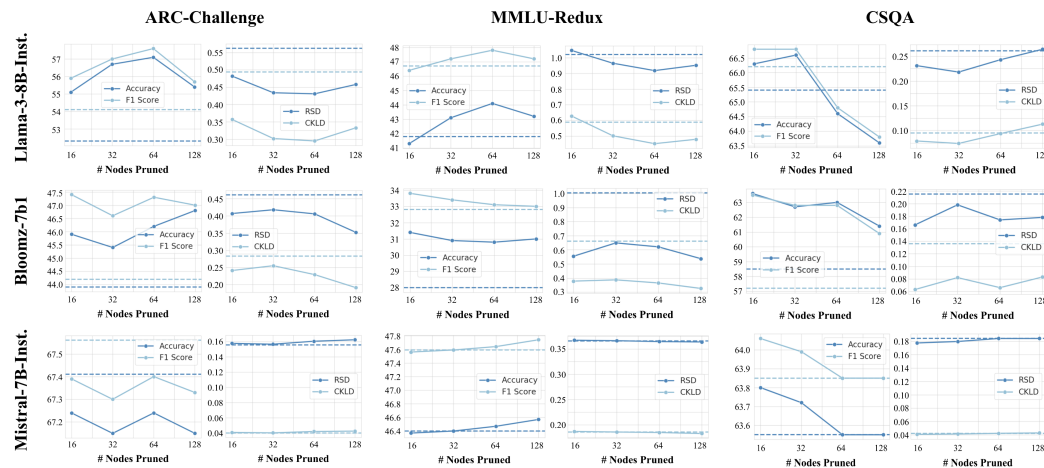


Figure 9: **Full list of plots on the number of nodes pruned.**

Table 8: **Different AOI setups.** The content, location, and number of auxiliary options are varied to see its effect with ARC-Challenge (top table), MMLU-Redux (middle table), and CSQA (bottom table).

Method	Llama-3-8B-Inst.				Bloomz-7b1				Mistral-7B-Inst.			
	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓
Model	52.3	54.1	0.562	0.494	43.9	44.2	0.461	0.283	67.4	67.6	0.156	0.040
Model + Ours	65.3	65.1	0.262	0.124	48.8	48.9	0.208	0.088	69.5	69.5	0.108	0.019
Arbitrary AOI	63.4	61.2	0.572	0.179	50.1	50.2	0.548	0.077	11.4	3.9	1.008	2.075
2-Choices AOI	70.2	69.9	0.175	0.067	46.3	47.6	0.381	0.198	69.0	69.0	0.131	0.031
3-Choices AOI	71.9	71.7	0.130	0.039	45.1	46.6	0.418	0.243	68.3	68.3	0.140	0.038
4-Choices AOI	72.4	72.3	0.130	0.036	43.9	45.6	0.438	0.266	68.4	68.4	0.138	0.036
First Choice AOI	67.9	67.6	0.222	0.106	44.2	45.3	0.455	0.232	68.1	68.1	0.109	0.025

Method	Llama-3-8B-Inst.				Bloomz-7b1				Mistral-7B-Inst.			
	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓
Base Model	41.8	46.7	1.021	0.589	28.0	32.8	1.003	0.661	46.4	47.6	0.366	0.186
Base Model + Ours	48.3	50.5	0.531	0.288	32.0	33.3	0.672	0.205	48.6	49.3	0.309	0.140
Arbitrary AOI	45.6	46.5	0.790	0.366	28.0	26.1	0.618	0.314	9.7	3.9	0.762	1.888
2-Choices AOI	49.4	50.9	0.442	0.201	30.5	32.7	0.774	0.332	47.7	48.4	0.327	0.157
3-Choices AOI	50.6	51.8	0.387	0.151	30.4	33.4	0.838	0.435	47.5	48.0	0.317	0.159
4-Choices AOI	51.7	52.8	0.352	0.117	30.0	33.4	0.633	0.479	47.1	47.7	0.328	0.169
First Choice AOI	46.1	47.6	0.515	0.295	31.8	35.4	0.647	0.338	44.7	45.0	0.291	0.160

Method	Llama-3-8B-Inst.				Bloomz-7b1				Mistral-7B-Inst.			
	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓	Acc. ↑	F1 ↑	RSD ↓	CKLD ↓
Base Model	65.4	66.2	0.261	0.095	58.5	57.2	0.215	0.136	63.6	63.9	0.184	0.042
Base Model + Ours	68.1	68.2	0.174	0.049	64.9	64.9	0.159	0.052	66.8	66.8	0.099	0.016
Arbitrary AOI	67.9	68.0	0.486	0.049	67.6	67.5	0.144	0.043	5.1	0.9	0.851	2.854
2-Choices AOI	68.1	68.2	0.149	0.031	59.5	59.8	0.261	0.129	65.6	65.6	0.134	0.034
3-Choices AOI	70.0	70.3	0.150	0.028	59.4	59.9	0.273	0.132	65.3	65.2	0.123	0.033
4-Choices AOI	70.4	70.5	0.137	0.023	58.7	59.4	0.282	0.130	64.8	64.7	0.137	0.038
First Choice AOI	69.5	69.4	0.142	0.037	48.5	52.7	0.602	0.713	66.2	66.3	0.118	0.018

**Location of the auxiliary option does not decide performance.** The location of the auxiliary option is another factor to consider. In our main experiments, we have appended the “I don’t know” option to the end of the choice list. In comparison, we try placing it in the first choice option (*i.e.*, with choice symbol ‘A’), corresponding to ‘First Choice AOI’ in Table 4. Overall, there were mixed results, indicating that the location of the auxiliary option is not a decisive factor in determining performance.

## F QUALITATIVE EXAMPLES

Here, we provide more qualitative examples to show how model response changes when our methods are applied. The examples are retrieved using the Llama-3-8B-Instruct model on the ARC-Challenge dataset. As observed in Figure 3(a), the original Llama-3 response is skewed towards ‘D’. The provided examples align with the result, and such ungrounded preference is debiased via our BNP+AOI.

**Original Question:** 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 gravity will become stronger. (D) Planetary days will become shorter.

⇒ **Base Model Response:** (D)

**Ground-truth:** (D)

**Permuted Question:** 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.

1134	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1135		
1136		
1137	<b>Original Question:</b> Petrified palm trees are found in sedimentary rock near glaciers. The	
1138	presence of the petrified palm trees most likely provides evidence for which statement? (A)	
1139	There was once more water in the area. (B) The area was once grassland. (C) There are active	
1140	faults in the area. (D) The climate in the area was once tropical.	
1141	⇒ <b>Base Model Response:</b> (D)	<b>Ground-truth:</b> (D)
1142	<b>Permuted Question:</b> Petrified palm trees are found in sedimentary rock near glaciers. The	
1143	presence of the petrified palm trees most likely provides evidence for which statement? (A)	
1144	There was once more water in the area. (B) The area was once grassland. (C) The climate in	
1145	the area was once tropical. (D) There are active faults in the area.	
1146	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1147		
1148		
1149	<b>Original Question:</b> According to cell classification, prokaryotic cells are separated from eu-	
1150	karyotic cells. Which feature is often used to distinguish prokaryotic cells from eukaryotic	
1151	cells? (A) plasma membranes (B) size differences (C) life processes (D) energy molecules	
1152	⇒ <b>Base Model Response:</b> (B)	<b>Ground-truth:</b> (B)
1153	<b>Permuted Question:</b> According to cell classification, prokaryotic cells are separated from	
1154	eukaryotic cells. Which feature is often used to distinguish prokaryotic cells from eukaryotic	
1155	cells? (A) life processes (B) size differences (C) plasma membranes (D) energy molecules	
1156	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (B)	<b>Ground-truth:</b> (B)
1157		
1158		
1159	<b>Original Question:</b> The morning temperature in a city is 41°F. If a sunny, mild day is forecast,	
1160	which temperature is most likely for 2:00 p.m.? (A) 32° F (B) 78° F (C) 98° F (D) 41° F	
1161	⇒ <b>Base Model Response:</b> (B)	<b>Ground-truth:</b> (B)
1162	<b>Permuted Question:</b> The morning temperature in a city is 41°F. If a sunny, mild day is fore-	
1163	cast, which temperature is most likely for 2:00 p.m.? (A) 32° F (B) 41° F (C) 78° F (D) 98°	
1164	F	
1165	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1166		
1167		
1168	<b>Original Question:</b> All natural resources on Earth are either renewable or nonrenewable.	
1169	Whether a resource is renewable or nonrenewable depends on how fast or slow the resource	
1170	is replaced. If the resource is used faster than it is replaced, then the resource will, in time,	
1171	disappear. Which activity shows the use of a nonrenewable natural resource? (A) A group of	
1172	people swims in a river. (B) A person bakes a cake with electricity produced by a hydroelectric	
1173	power plant. (C) A farmer grows vegetables to sell at a local market. (D) A construction crew	
1174	builds an iron bridge.	
1175	⇒ <b>Base Model Response:</b> (D)	<b>Ground-truth:</b> (D)
1176	<b>Permuted Question:</b> All natural resources on Earth are either renewable or nonrenewable.	
1177	Whether a resource is renewable or nonrenewable depends on how fast or slow the resource	
1178	is replaced. If the resource is used faster than it is replaced, then the resource will, in time,	
1179	disappear. Which activity shows the use of a nonrenewable natural resource? (A) A group of	
1180	people swims in a river. (B) A construction crew builds an iron bridge. (C) A farmer grows	
1181	vegetables to sell at a local market. (D) A person bakes a cake with electricity produced by a	
1182	hydroelectric power plant.	
1183	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (B)	<b>Ground-truth:</b> (B)
1184		
1185		
1186	<b>Original Question:</b> At which temperature does water freeze? (A) 32 degrees Celsius (B) 0	
1187	degrees Celsius (C) 100 degrees Celsius (D) 212 degrees Celsius	

1188	⇒ <b>Base Model Response:</b> (B)	<b>Ground-truth:</b> (B)
1189		
1190	<b>Permuted Question:</b> At which temperature does water freeze? (A) 0 degrees Celsius (B) 32	
1191	degrees Celsius (C) 100 degrees Celsius (D) 212 degrees Celsius	
1192	⇒ <b>Base Model Response:</b> (B) / <b>BNP+AOI Response:</b> (A)	<b>Ground-truth:</b> (A)
1193		
1194	<b>Original Question:</b> Fossil bones and teeth of dinosaurs have been researched for the last	
1195	century. Recent discoveries of fossilized dinosaurs have also revealed details of soft tissues,	
1196	such as skin. Which is best for a scientist to do when reporting research on dinosaurs now? (A)	
1197	exclude research on teeth or bones (B) delete earlier reports that were missing the new findings	
1198	(C) predict what the next discovery will be (D) analyze new data as it becomes available	
1199	⇒ <b>Base Model Response:</b> (D)	<b>Ground-truth:</b> (D)
1200		
1201	<b>Permuted Question:</b> Fossil bones and teeth of dinosaurs have been researched for the last	
1202	century. Recent discoveries of fossilized dinosaurs have also revealed details of soft tissues,	
1203	such as skin. Which is best for a scientist to do when reporting research on dinosaurs now?	
1204	(A) exclude research on teeth or bones (B) predict what the next discovery will be (C) analyze	
1205	new data as it becomes available (D) delete earlier reports that were missing the new findings	
1206	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1207		
1208	<b>Original Question:</b> What is the main function of photosynthetic cells within a plant? (A) to	
1209	change oxygen into carbon dioxide (B) to allow the passage of carbon dioxide into the plant (C)	
1210	to convert energy from sunlight into food energy (D) to break down sugar into usable chemicals	
1211	⇒ <b>Base Model Response:</b> (C)	<b>Ground-truth:</b> (C)
1212		
1213	<b>Permuted Question:</b> What is the main function of photosynthetic cells within a plant? (A)	
1214	to change oxygen into carbon dioxide (B) to break down sugar into usable chemicals (C) to	
1215	convert energy from sunlight into food energy (D) to allow the passage of carbon dioxide into	
1216	the plant	
1217	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1218		
1219	<b>Original Question:</b> What is the mass of a carbon atom that has 6 protons, 7 neutrons, and 6	
1220	electrons? (A) 7 (B) 19 (C) 6 (D) 13	
1221	⇒ <b>Base Model Response:</b> (D)	<b>Ground-truth:</b> (D)
1222		
1223	<b>Permuted Question:</b> What is the mass of a carbon atom that has 6 protons, 7 neutrons, and 6	
1224	electrons? (A) 6 (B) 7 (C) 13 (D) 19	
1225	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1226		
1227	<b>Original Question:</b> Air has no color and cannot be seen, yet it takes up space. What could	
1228	be done to show that air takes up space? (A) observe clouds forming (B) blow up a beach ball	
1229	or balloon (C) measure the air temperature (D) weigh a glass before and after it is filled with	
1230	water	
1231	⇒ <b>Base Model Response:</b> (B)	<b>Ground-truth:</b> (B)
1232		
1233	<b>Permuted Question:</b> Air has no color and cannot be seen, yet it takes up space. What could	
1234	be done to show that air takes up space? (A) observe clouds forming (B) measure the air	
1235	temperature (C) blow up a beach ball or balloon (D) weigh a glass before and after it is filled	
1236	with water	
1237	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1238		
1239	<b>Original Question:</b> Which geologic process most likely caused the formation of the Mount	
1240	St. Helens Volcano? (A) diverging boundaries (B) converging boundaries (C) transform faults	
1241	(D) rift zone	

1242 ⇒ **Base Model Response:** (B) **Ground-truth:** (B)  
 1243  
 1244 **Permuted Question:** Which geologic process most likely caused the formation of the Mount  
 1245 St. Helens Volcano? (A) converging boundaries (B) diverging boundaries (C) transform faults  
 1246 (D) rift zones  
 1247 ⇒ **Base Model Response:** (D) / **BNP+AOI Response:** (A) **Ground-truth:** (A)  
 1248

1249 We also provide results with Bloomz-7b1 on ARC-Challenge. Similar to the trend shown in Figure 3  
 1250 (a), the original response is biased towards ‘A’, which is corrected through our debiasing approach.  
 1251

1252 **Original Question:** Devil facial tumor disease (DFTD) is a disease that is decimating the  
 1253 population of Tasmanian devils. The disease passes from one animal to another through bites  
 1254 and is caused by parasites. The parasites cause cancerous tumors that spread throughout an  
 1255 infected animal’s body and kill it. What is the best description of DFTD? (A) a non-infectious,  
 1256 cell-cycle disease (B) a non-infectious, chronic disease (C) an infectious, cell-cycle disease (D)  
 1257 an infectious, chronic disease  
 1258 ⇒ **Base Model Response:** (C) **Ground-truth:** (C)  
 1259  
 1260 **Permuted Question:** Devil facial tumor disease (DFTD) is a disease that is decimating the  
 1261 population of Tasmanian devils. The disease passes from one animal to another through bites  
 1262 and is caused by parasites. The parasites cause cancerous tumors that spread throughout an  
 1263 infected animal’s body and kill it. What is the best description of DFTD? (A) a non-infectious,  
 1264 cell-cycle disease (B) an infectious, cell-cycle disease (C) a non-infectious, chronic disease (D)  
 1265 an infectious, chronic disease  
 1266 ⇒ **Base Model Response:** (A) / **BNP+AOI Response:** (B) **Ground-truth:** (B)  
 1267

1268 **Original Question:** Which of these gases is the most abundant greenhouse gas in the lower  
 1269 atmosphere of Earth? (A) carbon dioxide (B) methane (C) water vapor (D) ozone  
 1270 ⇒ **Base Model Response:** (C) **Ground-truth:** (C)  
 1271  
 1272 **Permuted Question:** Which of these gases is the most abundant greenhouse gas in the lower  
 1273 atmosphere of Earth? (A) ozone (B) methane (C) water vapor (D) carbon dioxide  
 1274 ⇒ **Base Model Response:** (D) / **BNP+AOI Response:** (C) **Ground-truth:** (C)  
 1275

1276 **Original Question:** It was once thought that living organisms could come from non-living  
 1277 matter. For example, people believed that flies would develop from rotting meat. This idea was  
 1278 later disproved primarily because of (A) the discovery of the atom. (B) continued experimen-  
 1279 tation. (C) better surgical techniques. (D) the invention of the microscope.  
 1280 ⇒ **Base Model Response:** (B) **Ground-truth:** (B)  
 1281  
 1282 **Permuted Question:** It was once thought that living organisms could come from non-living  
 1283 matter. For example, people believed that flies would develop from rotting meat. This idea  
 1284 was later disproved primarily because of (A) the discovery of the atom. (B) better surgical  
 1285 techniques. (C) continued experimentation. (D) the invention of the microscope  
 1286 ⇒ **Base Model Response:** (A) / **BNP+AOI Response:** (C) **Ground-truth:** (C)  
 1287

1288 **Original Question:** In the spring and early summer, bears often scratch their backs against  
 1289 trees to remove winter fur. This is an example of an animal (A) responding to its environment  
 1290 (B) beginning hibernation (C) completing its life cycle (D) preparing for migration  
 1291 ⇒ **Base Model Response:** (A) **Ground-truth:** (A)  
 1292  
 1293 **Permuted Question:** In the spring and early summer, bears often scratch their backs against  
 1294 trees to remove winter fur. This is an example of an animal (A) completing its life cycle (B)  
 1295 beginning hibernation (C) responding to its environment (D) preparing for migration



1296	⇒ <b>Base Model Response:</b> (A) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1297		
1298	<b>Original Question:</b> Which tool would be best to use to determine how long it takes a cup of	
1299	water to boil? (A) balance (B) hot plate (C) thermometer (D) stopwatch	
1300		
1301	⇒ <b>Base Model Response:</b> (D)	<b>Ground-truth:</b> (D)
1302	<b>Permuted Question:</b> Which tool would be best to use to determine how long it takes a cup of	
1303	water to boil? (A) balance (B) hot plate (C) stopwatch (D) thermometer	
1304		
1305	⇒ <b>Base Model Response:</b> (D) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1306		
1307	<b>Original Question:</b> The salt in ocean water comes from all of the following except (A) melting	
1308	glacial ice. (B) volcanic emissions. (C) eroding land. (D) reactions on the sea floor.	
1309		
1310	⇒ <b>Base Model Response:</b> (A)	<b>Ground-truth:</b> (A)
1311	<b>Permuted Question:</b> The salt in ocean water comes from all of the following except (A)	
1312	eroding land. (B) melting glacial ice. (C) volcanic emissions. (D) reactions on the sea floor.	
1313		
1314	⇒ <b>Base Model Response:</b> (A) / <b>BNP+AOI Response:</b> (B)	<b>Ground-truth:</b> (B)
1315		
1316	<b>Original Question:</b> Which is most useful to a student who is separating aluminum screws	
1317	from steel screws? (A) a screen filter (B) a large funnel (C) a magnifying glass (D) a horseshoe	
1318	magnet	
1319	⇒ <b>Base Model Response:</b> (D)	<b>Ground-truth:</b> (D)
1320	<b>Permuted Question:</b> Which is most useful to a student who is separating aluminum screws	
1321	from steel screws? (A) a large funnel (B) a screen filter (C) a horseshoe magnet (D) a magni-	
1322	fying glass	
1323		
1324	⇒ <b>Base Model Response:</b> (A) / <b>BNP+AOI Response:</b> (C)	<b>Ground-truth:</b> (C)
1325		
1326	<b>Original Question:</b> Over a long period of time, running water in a river erodes the riverbed.	
1327	This erosion causes the river to (A) move faster and cleaner. (B) become deeper and wider. (C)	
1328	stop flowing. (D) create waves	
1329	⇒ <b>Base Model Response:</b> (B)	<b>Ground-truth:</b> (B)
1330	<b>Permuted Question:</b> Over a long period of time, running water in a river erodes the riverbed.	
1331	This erosion causes the river to (A) stop flowing. (B) create waves. (C) move faster and cleaner.	
1332	(D) become deeper and wider.	
1333		
1334	⇒ <b>Base Model Response:</b> (A) / <b>BNP+AOI Response:</b> (D)	<b>Ground-truth:</b> (D)
1335		
1336	<b>Original Question:</b> A student examined diagrams of two different cells. One cell was prokary-	
1337	otic, and the other cell was eukaryotic. What should the student do to identify a major differ-	
1338	ence between the diagrams? (A) check to see which diagram shows a nucleus (B) check to	
1339	see which diagram shows cytoplasm (C) compare the shapes of the two cells (D) compare the	
1340	number of vacuoles in the two cells	
1341	⇒ <b>Base Model Response:</b> (A)	<b>Ground-truth:</b> (A)
1342	<b>Permuted Question:</b> A student examined diagrams of two different cells. One cell was	
1343	prokaryotic, and the other cell was eukaryotic. What should the student do to identify a major	
1344	difference between the diagrams? (A) compare the shapes of the two cells (B) check to see	
1345	which diagram shows a nucleus (C) check to see which diagram shows cytoplasm (D) compare	
1346	the number of vacuoles in the two cells	
1347		
1348	⇒ <b>Base Model Response:</b> (A) / <b>BNP+AOI Response:</b> (B)	<b>Ground-truth:</b> (B)
1349		

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

**Original Question:** Which structures are common to both plant and animal cells? (A) cell membrane, nucleus, mitochondrion (B) vacuole, chloroplast, nucleus (C) nucleus, cell wall, cell membrane (D) mitochondrion, vacuole, cell wall

⇒ **Base Model Response:** (A)

**Ground-truth:** (A)

**Permuted Question:** Which structures are common to both plant and animal cells? (A) vacuole, chloroplast, nucleus (B) cell membrane, nucleus, mitochondrion (C) nucleus, cell wall, cell membrane (D) mitochondrion, vacuole, cell wall

⇒ **Base Model Response:** (A) / **BNP+AOI Response:** (B)

**Ground-truth:** (B)

**Original Question:** Students use tweezers and magnifying glasses to examine a piece of mold on bread. Which should they also use for safety in this investigation? (A) bright light (B) breathing masks (C) dark glasses (D) hot plates

⇒ **Base Model Response:** (B)

**Ground-truth:** (B)

**Permuted Question:** Students use tweezers and magnifying glasses to examine a piece of mold on bread. Which should they also use for safety in this investigation? (A) bright light (B) dark glasses (C) breathing masks (D) hot plates

⇒ **Base Model Response:** (A) / **BNP+AOI Response:** (C)

**Ground-truth:** (C)

**Original Question:** In 1903 Mary Anderson invented the first windshield wiper. How did this invention most likely help people? (A) It made cars easier for people to buy. (B) It kept people from driving too fast. (C) It helped people use less gas. (D) It made cars safer to drive in bad weather.

⇒ **Base Model Response:** (D)

**Ground-truth:** (D)

**Permuted Question:** In 1903 Mary Anderson invented the first windshield wiper. How did this invention most likely help people? (A) It helped people use less gas. (B) It kept people from driving too fast. (C) It made cars easier for people to buy. (D) It made cars safer to drive in bad weather.

⇒ **Base Model Response:** (A) / **BNP+AOI Response:** (D)

**Ground-truth:** (D)

## G LIMITATIONS AND BROADER IMPACT

**Limitations.** The limitation of this work (and also most works on mitigating selection bias) is that we still do not know the root cause of the selection bias. While there have been various hypotheses on the reason behind this phenomenon, most focused on the superficial effect of it without considering *what in the first place triggered such ungrounded preferences*. Future research will need to unravel the core of selection bias by answering questions like, *What data points cause selection bias?* or *What makes the difference in choice preferences between heterogenous model families?* These questions will be critical in understanding LLMs in general, as it is closely related to how the models choose the next tokens to output.

**Broader Impact.** This work reveals and mitigates a type of bias present in recent large language models (LLMs). Considering that LLMs have become an integral part of various applications from customer service to science, the presence of any type of bias can negatively impact the reliability of systems and degrade precision in model- or data-driven decision-making. By addressing the bias, our research not only improves the accuracy and fairness of these models but also has the potential to enhance the trustworthiness of LLMs in general. Moreover, this work serves as a foundation for ongoing efforts to scrutinize and enhance LLM-automated systems, introducing a new perspective on analyzing performance.

**Future Application.** One closely related application of selection bias debiasing is data annotation. Many works discussed ways to leverage LLMs for automated annotation (He et al., 2024; Eckman

1404 et al., 2024), or devised human-machine collaborative frameworks (Li et al., 2023a). We expect our  
1405 work to benefit such annotation systems by reducing the selection bias in answering MCQs.  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457