

UniBias: Unveiling and Mitigating LLM Bias through Internal Attention and FFN Manipulation

Hanzhang Zhou, Zijian Feng,
Zixiao Zhu, Junlang Qian,
Kezhi Mao



**NANYANG
TECHNOLOGICAL
UNIVERSITY**

SINGAPORE

ETH zürich
(FRS) FUTURE
RESILIENT
SYSTEMS

Overview

Problem: Prompt brittleness & LLM bias

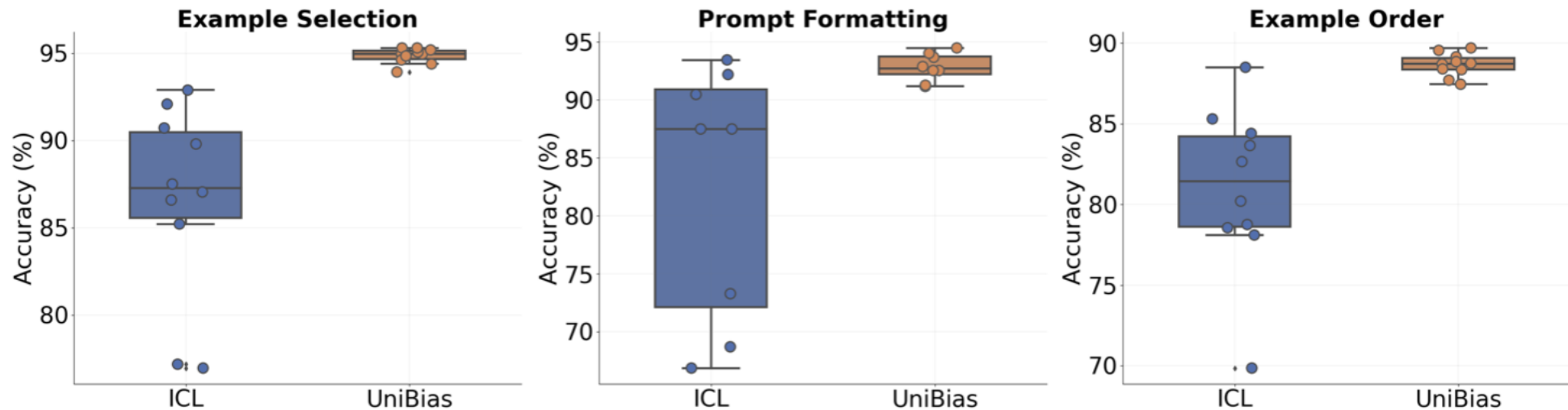
Highlight:

- Unveiling the internal mechanism of LLM bias
- Mitigating LLM bias by manipulation of attention heads and FFN vectors, offering an alternative to calibration methods that rely on external adjustment of LLM outputs.
- Demonstrating an effective way to steer LLM behavior through manipulation of LLM internal structures

Problem Definition

Prompt Brittleness: LLMs are highly sensitive to the choice and order of examples, and prompt formatting, undermining the robustness and adaptability of LLMs

LLM Bias: Prompt brittleness is found to arise from the bias in LLMs toward predicting certain answers



Unibias – Interpret the Contribution of LLM Components

How to interpret the contribution of LLM internal components to LLM prediction?

Mechanistic Interpretability

Quantify the contribution of FFN vector and attention heads to prediction

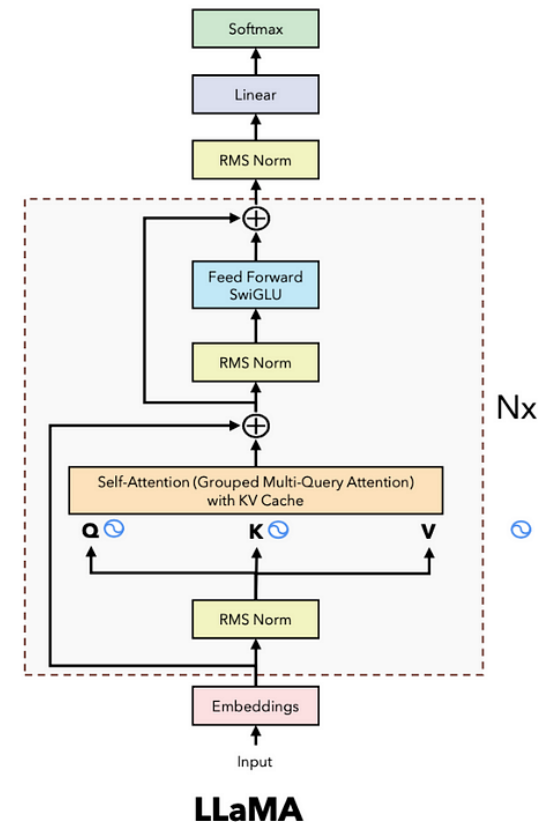
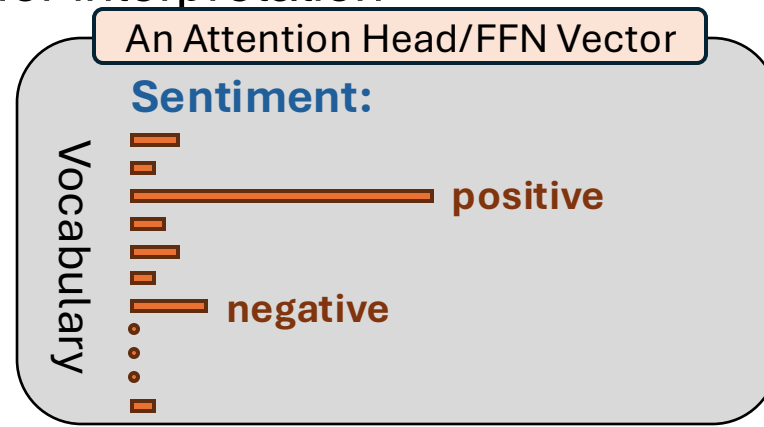
- **Attention Heads:** The output of an attention layer can be viewed as the sum of its respective attention heads

$$\text{Att}^\ell(X^\ell) = \text{Concat} \left[A^{\ell,1} X^\ell W_V^{\ell,1}, A^{\ell,2} X^\ell W_V^{\ell,2}, \dots, A^{\ell,H} X^\ell W_V^{\ell,H} \right] W_O^\ell = \sum_{j=1}^H A^{\ell,j} (X^\ell W_V^{\ell,j}) W_O^{\ell,j}$$

- **FFN Vectors:** The output of an FFN layer can be viewed as the weighted sum of its FFN value vectors

$$\text{FFN}^\ell(\mathbf{x}^\ell) = f(\mathbf{x}^\ell \mathbf{K}^{\ell T}) \mathbf{V}^\ell = \sum_{i=1}^{d_m} f(\mathbf{x}^\ell \cdot \mathbf{k}_i^\ell) \mathbf{v}_i^\ell = \sum_{i=1}^{d_m} m_i^\ell \mathbf{v}_i^\ell$$

- **Logit Lens:** Map hidden states into the vocabulary space using the unembedding matrix for interpretation

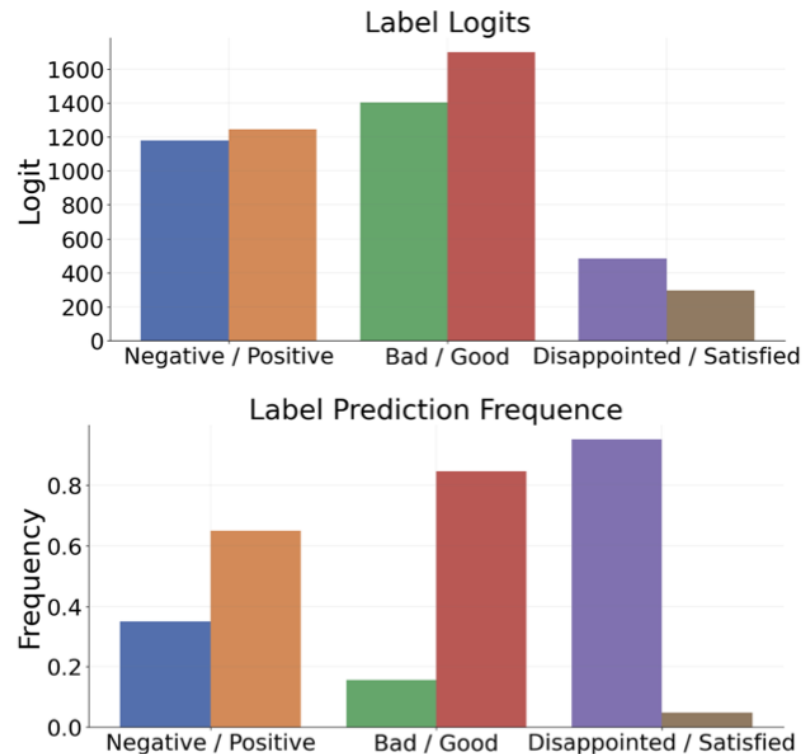


Unibias – Internal Mechanism of LLM bias

Vanilla Label Bias: Uncontextual preference of the model towards predicting certain label names

Internal Mechanism: A corresponding uncontextual preference is identified on FFN layers

- Accumulate the logits for all FFN value vectors
- The accumulated logits are uncontextual and biased toward certain label names

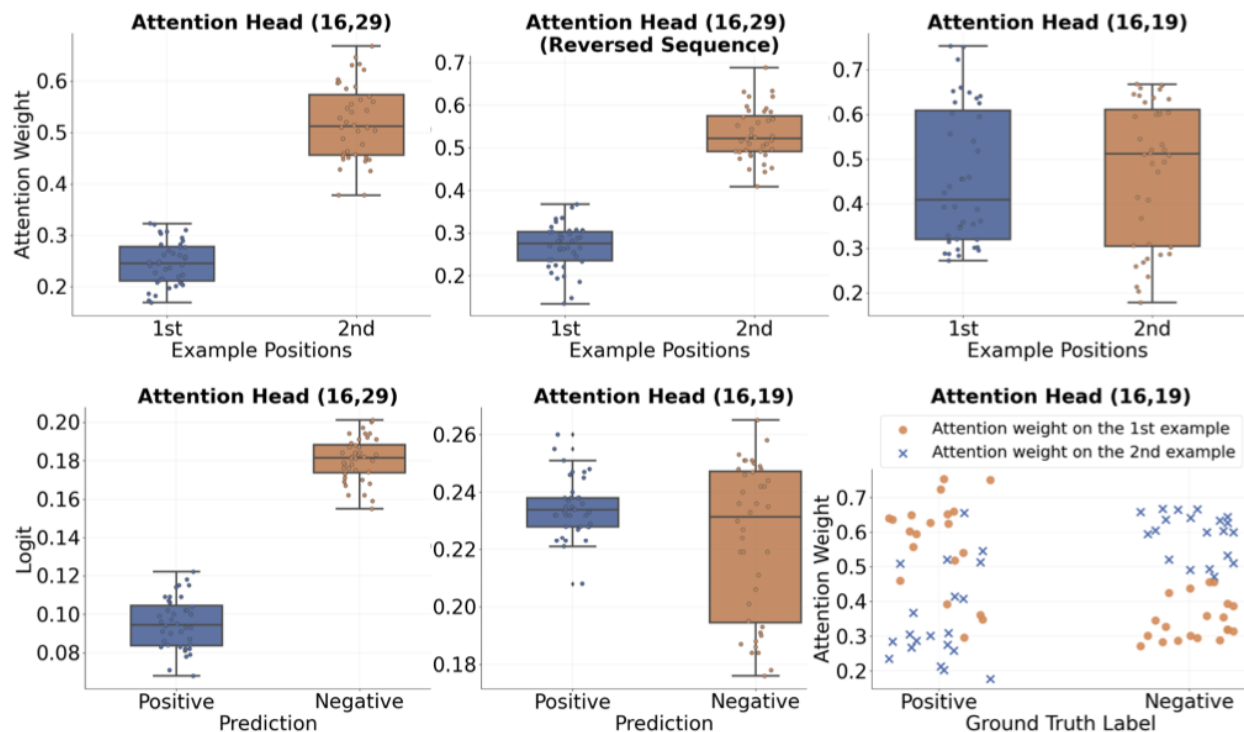


Internal Mechanism of LLM bias

Position Bias: LLMs tend to favor the label of the example at the end of the prompt

Internal Mechanism: Specific attention heads consistently prioritize the example at the end of the prompt

- Comparison between a biased attention head (layer 16, head 29) and an unbiased attention head (layer 16, head 19)



Methodology

Findings: Various bias factors stem from the biased behaviors of attention heads and FFN vectors.

**Can we identify the biased components of LLMs
and mitigate their impact on label prediction?**

Methodology

How to identify an FFN vector or an attention head as biased?

It consistently introduces a biased preference towards specific labels into the residual stream, regardless of variations in the test samples.

- **Relatedness Criterion:** The information introduced by the FFN vector (or attention head) should closely relate to label prediction.
- **Biased Criterion:** The information contributed to the residual stream by the FFN vector (or attention head) exhibits a biased distribution, favoring certain labels over others.
- **Low Variance Criterion:** The label prediction information added by the FFN vector (or attention head) to the residual stream is almost identical across test samples with different labels,

$$\left\{ \begin{array}{l} \frac{1}{p} \sum_{k=0}^{p-1} \text{Sum}(\mathbf{G}_{k,:}) = \frac{1}{p} \sum_{k=0}^{p-1} \text{Sum}(\mathbf{g}^{(k)}) = \frac{1}{p} \sum_{k=0}^{p-1} \sum_{j=0}^{c-1} g_j^{(k)} > th_{FFN}^1 \\ \frac{1}{p} \sum_{k=0}^{p-1} \text{Bias}(\mathbf{G}_{k,:}) = \frac{1}{p} \sum_{k=0}^{p-1} \text{Bias}(\mathbf{g}^{(k)}) = \frac{1}{p} \frac{1}{c} \sum_{k=0}^{p-1} \sum_{j=0}^{c-1} (g_j^{(k)} - \mu(\mathbf{g}^{(k)})) > th_{FFN}^2 \\ CV(\mathbf{m}) = \frac{\sigma(\mathbf{m})}{\mu(\mathbf{m})} = \frac{\sqrt{\frac{1}{p} \sum_{j=0}^{p-1} (m_k - \mu(\mathbf{m}))^2}}{\frac{1}{p} \sum_{k=0}^{p-1} m_k} < th_{FFN}^3 \end{array} \right. \quad \left\{ \begin{array}{l} \frac{1}{p} \sum_{k=0}^{p-1} \text{Sum}(A_{k,:}) = \frac{1}{p} \sum_{k=0}^{p-1} \text{Sum}(\mathbf{a}^{(k)}) = \frac{1}{p} \sum_{k=0}^{p-1} \sum_{j=1}^c a_j^{(k)} > th_{Att}^1 \\ \frac{1}{p} \sum_{k=0}^{p-1} \text{Bias}(A_{k,:}) = \frac{1}{p} \sum_{k=0}^{p-1} \text{Bias}(\mathbf{a}^{(k)}) = \frac{1}{p} \frac{1}{c} \sum_{k=0}^{p-1} \sum_{j=0}^{c-1} (a_j^{(k)} - \mu(\mathbf{a}^{(k)})) > th_{Att}^2 \\ \sum_{j=0}^{c-1} w_j \cdot CV(A_{:,j}) = w_j \cdot \frac{\sigma(A_{:,j})}{\mu(A_{:,j})} < th_{Att}^3 \end{array} \right.$$

How to Eliminate the influence of biased FFN vectors and attention heads? Masking

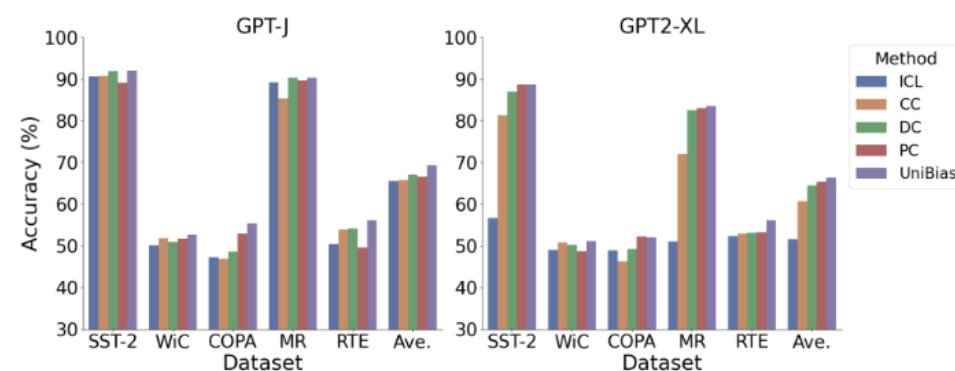
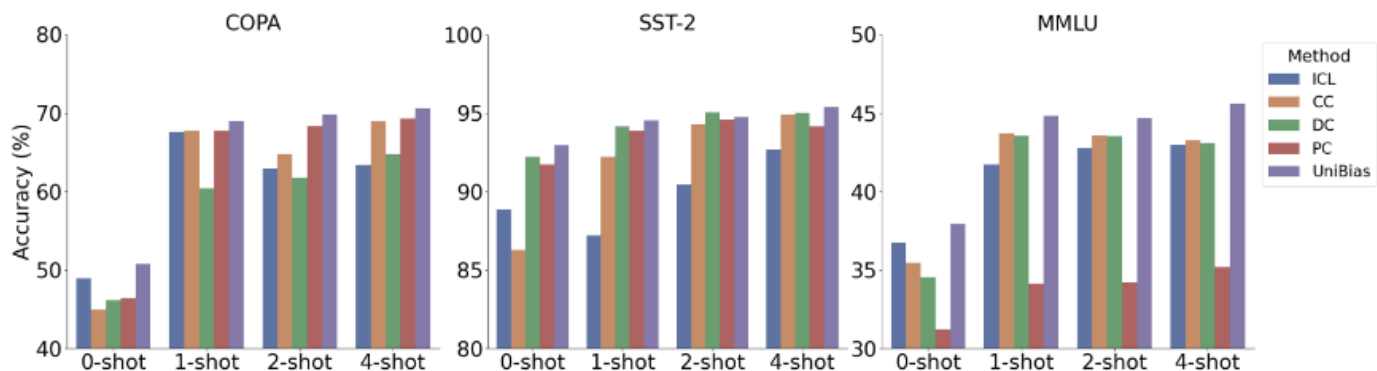
Experiments

Dataset: 12 datasets (reasoning, sentiment analysis, topic classification, natural language inference, and word disambiguation)

LLMs: 4 LLMs

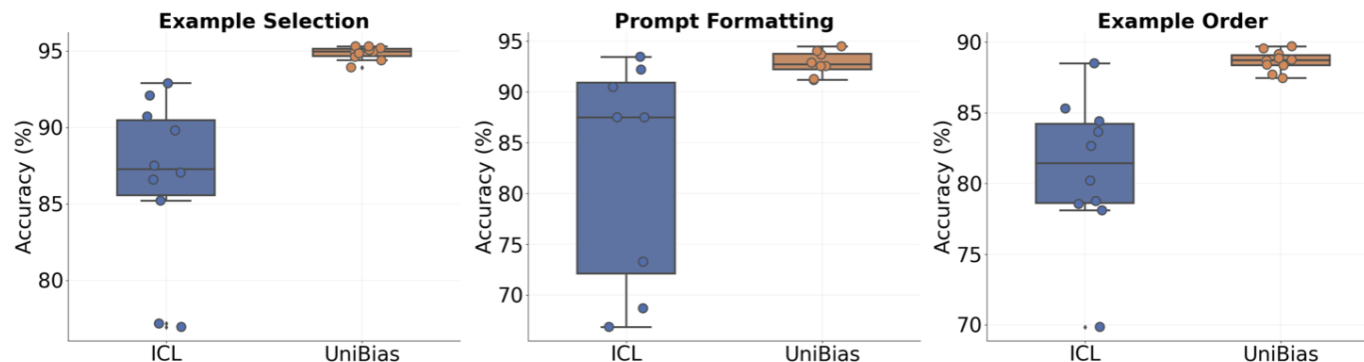
Baselines: Calibration methods (CC, DC, PC)

Dataset	Llama-2 7b					Llama-2 13b				
	ICL	CC	DC	PC	UniBias	ICL	CC	DC	PC	UniBias
SST-2	87.22 _{6.03}	92.24 _{3.39}	94.15 _{1.22}	93.90 _{1.54}	94.54 _{0.62}	93.90 _{1.79}	95.25 _{0.93}	95.37 _{0.70}	94.56 _{1.71}	95.46 _{0.52}
MNLI	53.83 _{2.22}	53.36 _{3.16}	52.19 _{2.55}	45.38 _{5.01}	54.97 _{0.88}	62.43 _{1.49}	63.89 _{0.81}	61.86 _{1.23}	57.47 _{3.53}	64.65 _{2.73}
WiC	50.00 _{0.16}	52.19 _{2.00}	52.40 _{1.69}	57.11 _{2.49}	53.71 _{1.16}	54.48 _{3.19}	50.63 _{1.73}	49.72 _{0.30}	55.67 _{1.67}	57.93 _{1.70}
COPA	67.60 _{2.30}	67.80 _{2.17}	60.40 _{2.79}	67.80 _{3.70}	69.00 _{2.74}	67.50 _{10.40}	75.20 _{7.80}	71.00 _{8.80}	76.80 _{6.30}	83.20 _{2.70}
CR	91.54 _{0.39}	92.13 _{0.40}	92.61 _{0.44}	91.97 _{0.35}	92.61 _{0.11}	91.01 _{1.30}	92.13 _{0.88}	92.23 _{0.76}	91.65 _{0.64}	92.34 _{0.74}
AGNews	85.59 _{1.87}	83.54 _{1.96}	89.08 _{0.86}	86.81 _{2.92}	88.29 _{1.24}	89.14 _{0.44}	88.23 _{1.14}	89.34 _{0.61}	86.03 _{0.65}	88.68 _{0.43}
MR	89.37 _{1.83}	91.77 _{1.42}	92.35 _{0.23}	91.39 _{1.65}	92.19 _{0.37}	90.10 _{2.10}	93.20 _{0.57}	93.00 _{0.52}	92.80 _{0.86}	92.23 _{1.12}
RTE	66.21 _{7.30}	64.33 _{3.68}	65.49 _{2.09}	62.59 _{4.71}	67.65 _{6.44}	76.10 _{4.73}	71.99 _{5.02}	66.21 _{1.09}	75.31 _{2.90}	78.23 _{2.13}
SST-5	46.97 _{0.87}	51.36 _{1.69}	51.92 _{1.77}	55.41 _{1.51}	53.79 _{1.46}	51.03 _{1.25}	47.20 _{1.69}	48.98 _{2.11}	53.63 _{0.95}	51.80 _{1.00}
TREC	72.92 _{12.42}	76.44 _{3.21}	77.16 _{3.94}	74.92 _{5.78}	80.80 _{3.17}	74.70 _{12.10}	83.80 _{3.86}	80.50 _{9.07}	81.85 _{9.53}	81.25 _{6.86}
ARC	51.90 _{0.60}	53.10 _{0.40}	53.00 _{0.60}	40.40 _{0.50}	53.10 _{0.60}	66.54 _{0.33}	64.33 _{0.99}	64.88 _{0.59}	59.47 _{1.07}	66.81 _{0.37}
MMLU	41.73 _{2.25}	43.72 _{0.97}	43.57 _{1.38}	34.12 _{3.41}	44.83 _{0.24}	53.53 _{1.55}	50.84 _{1.57}	51.81 _{1.24}	45.50 _{1.65}	53.55 _{1.05}
Avg.	67.07	68.49	68.70	66.81	70.46	72.54	73.06	72.08	72.56	75.51

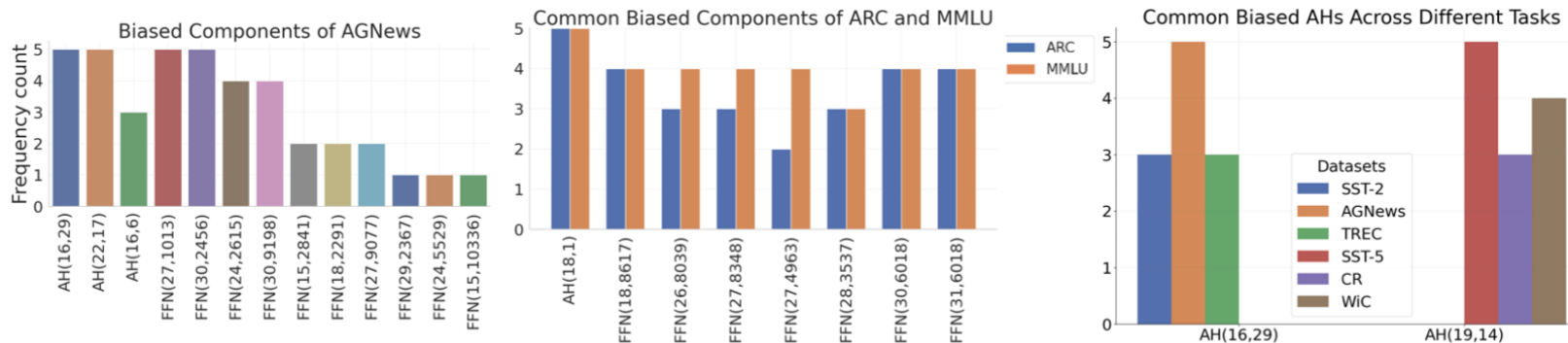


Experiments

Alleviating Prompt Brittleness:



Removing Common Biased Components



	SST2	MMLU	COPA	RTE	MR	Trec	Avg.
ICL	87.22 _{6.03}	41.73 _{2.25}	67.60 _{2.30}	66.21 _{7.30}	89.37 _{1.83}	72.92 _{12.42}	70.84
Unibias	94.54 _{0.62}	44.83 _{0.24}	69.00 _{2.74}	67.65 _{6.44}	92.19 _{0.37}	80.80 _{3.17}	74.84
Eliminating Common Biased Components	94.32 _{0.60}	44.20 _{1.14}	68.00 _{2.87}	67.37 _{4.60}	92.43 _{0.09}	77.60 _{4.75}	73.98

Summary

New Direction for LLM Bias Mitigation: Potentially stimulate future research on LLM bias mitigation from inner structure manipulation, offering a new direction for the field.

Unveil Internal Mechanisms of LLM Bias: Provide deep insights that go beyond superficial observations, revealing the inner causes of LLM bias.

Enhance the accuracy and robustness of ICL: Effectively address prompt brittleness, thereby boosting the robustness of LLMs and markedly improving the accuracy of ICL.

LLM manipulation: Showcase how to steer LLM behavior by manipulate the internal structures of large language models