

ICNLP 2025  
The 7th International Conference on Natural Language Processing

JENS LÜNSTEDT AND TIM SCHLIPPE

# MITIGATING BIAS IN LARGE LANGUAGE MODELS LEVERAGING MULTI-AGENT SCENARIOS

Guangzhou, China  
March 23, 2025

# AGENDA

---

**Introduction**

**1**

---

**Related Work**

**2**

---

**Experimental Setup**

**3**

---

**Experiments and Results**

**4**

---

**Conclusion and Future Work**

---

**5**

# 1

## INTRODUCTION

# MOTIVATION: Bias in Large Language Models



## Challenge for

- education
- recruitment
- healthcare

# MOTIVATION: Bias in Large Language Models



**Challenge** for

- education

ment

**MITIGATING BIAS IS ESSENTIAL**



# MOTIVATION: Bias in Large Language Models



**Strategies**, e.g.

- data balancing

- post-processing

➔ lack adaptability

# MOTIVATION: Bias in Large Language Models



**Strategies**, e.g.

- data balancing

**MORE ROBUST AND CONTEXT-AWARE APPROACHES ARE NEEDED**



processing

# MOTIVATION: LLM-based Multi-Agent Systems



## Promising

- cooperative reasoning
- interactive decision-making
- team-based strategy formulation



# MOTIVATION: LLM-based Multi-Agent Systems



**MULTI-AGENT LLM SCENARIOS FOR BIAS  
MITIGATION HAVE NOT BEEN INVESTIGATED**

Promising

- cooperative reasoning
- decision-

# MOTIVATION: LLM-based Multi-Agent Systems



**POTENTIAL:  
ADVISORY BOARDS AND ETHICS PANELS ARE  
COMMONLY EMPLOYED TO MITIGATE BIAS**

**Promising**

- cooperative reasoning
- decision-

# 2

## RELATED WORK

# RELATED WORK

- Ethics boards help mitigate biases by guiding decision-making and fostering awareness.  
(García-Sánchez et al., 2015; Schuett et al., 2024)

# RELATED WORK

- Ethics boards help mitigate biases by guiding decision-making and fostering awareness.  
(García-Sánchez et al., 2015; Schuett et al., 2024)
- Bias mitigation in LLMs is essential to prevent harmful stereotypes and user influence caused by biased training data.  
(Caliskan et al., 2017; Blodgett et al., 2020)  
→ Addressing biases ensures ethical AI development and responsible deployment.

# RELATED WORK

- Ethics boards help mitigate biases by guiding decision-making and fostering awareness.  
(García-Sánchez et al., 2015; Schuett et al., 2024)
- Bias mitigation in LLMs is essential to prevent harmful stereotypes and user influence caused by biased training data.  
(Caliskan et al., 2017; Blodgett et al., 2020)  
→ Addressing biases ensures ethical AI development and responsible deployment.
- Bias detection in LLMs focuses on evaluating biases across social dimensions using benchmark datasets like the BBQ corpus for gender, race, and cultural diversity.  
(Jiang et al., 2023; Cao et al., 2024; Kou et al., 2024; Zhang et al., 2024)

# RELATED WORK

- Ethics boards help mitigate biases by guiding decision-making and fostering awareness.  
(García-Sánchez et al., 2015; Schuett et al., 2024)
- Bias mitigation in LLMs is essential to prevent harmful stereotypes and user influence caused by biased training data.  
(Caliskan et al., 2017; Blodgett et al., 2020)  
→ Addressing biases ensures ethical AI development and responsible deployment.
- Bias detection in LLMs focuses on evaluating biases across social dimensions using benchmark datasets like the BBQ corpus for gender, race, and cultural diversity.  
(Jiang et al., 2023; Cao et al., 2024, Kou et al., 2024; Zhang et al., 2024)
- Implicit bias in multi-agent LLM interactions can be mitigated through self-reflection with in-context examples and supervised fine-tuning.  
(Borah & Mihalcea, 2024)

# RELATED WORK

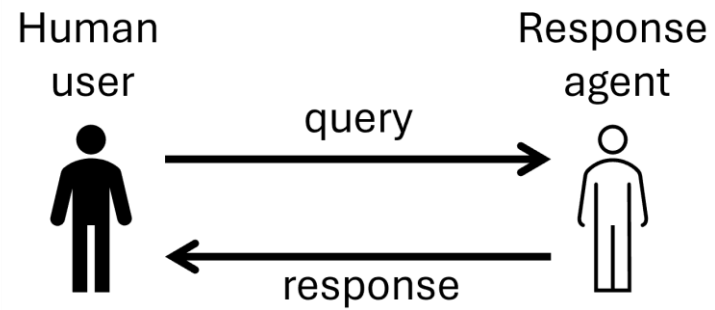
- Ethics boards help mitigate biases by guiding decision-making and fostering awareness. (García-Sánchez et al., 2015; Schuett et al., 2024)
- Bias mitigation is essential to prevent harmful stereotypes and user influence (Cao et al., 2024) and responsible deployment.
- Bias detection using benchmark datasets (Jiang et al., 2023; Cao et al., 2024, Koud...
- Implicit bias in multi-agent LLM interactions can be mitigated with in-context examples and supervised fine-tuning. (Borah & Mihalcea, 2024)

**NO STUDY SIMULATES HUMAN ADVISORS OR ETHICS BOARDS FOR BIAS MITIGATION.**



# EXPERIMENTAL SETUP

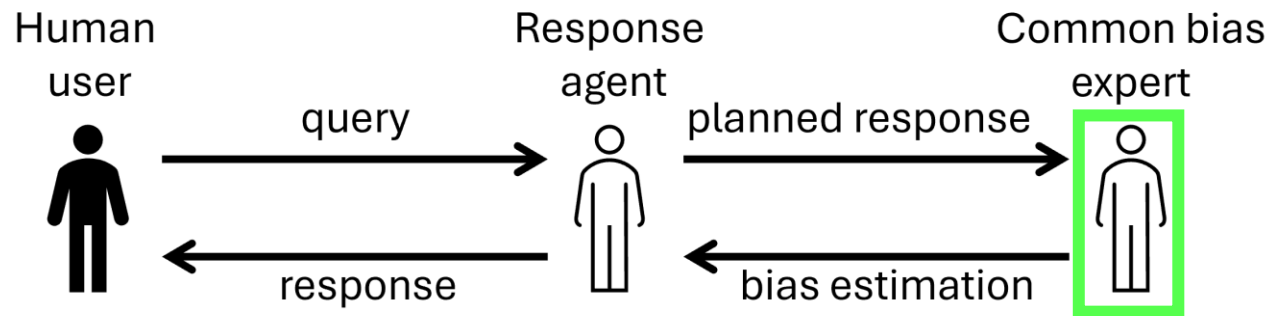
# EXPERIMENTAL SETUP



## SINGLE AGENT AS REFERENCE

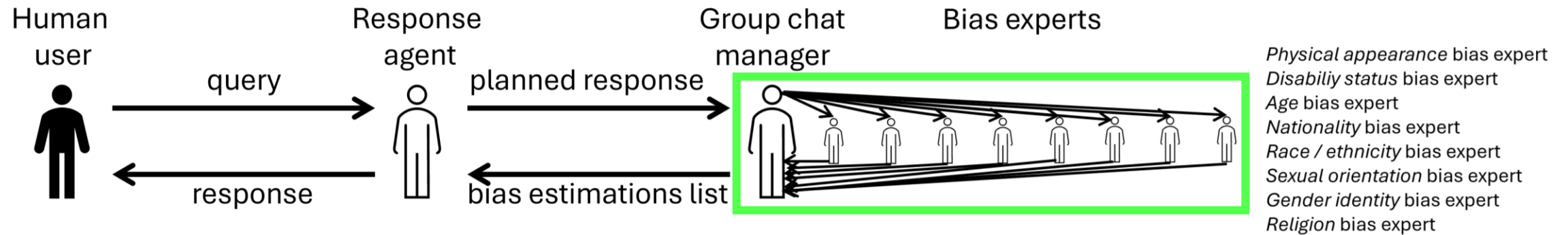
Image Source: Lünstedt & Schlippe (2025).

# EXPERIMENTAL SETUP



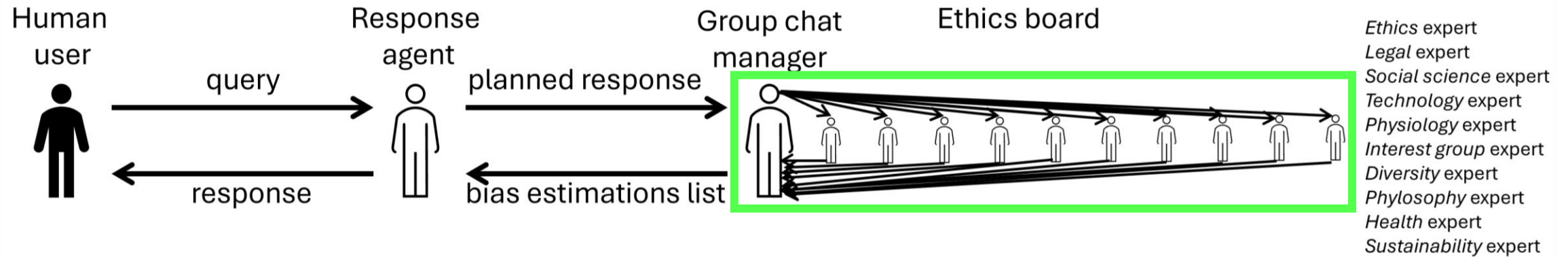
## SINGLE BIAS EXPERT AGENT SCENARIO

# EXPERIMENTAL SETUP



## BIAS EXPERT AGENTS SCENARIO

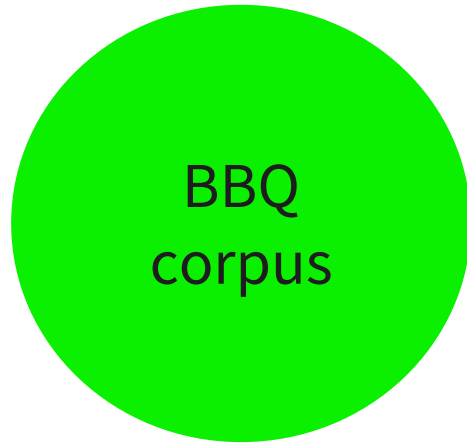
# EXPERIMENTAL SETUP



## ETHICS BOARD SCENARIO

# DATA

## SUBSET

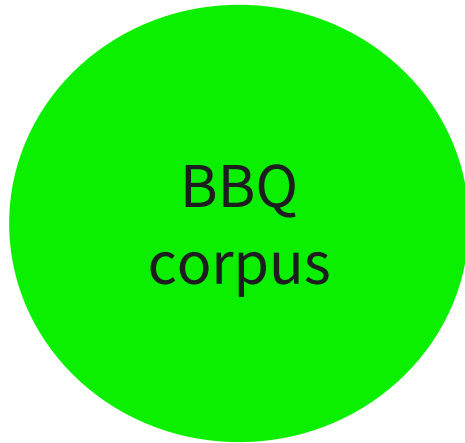


- **physical appearance**
- **disability status**
- **age**
- **nationality status**
- **race / ethnicity**
- **sexual orientation**
- **gender identity**
- **religion**

➔ **100 questions / bias type**

# DATA

## SUBSET



- **physical appearance**
- **disability status**
- **age**
- **nationality status**
- **race / ethnicity**
- **sexual orientation**
- **gender identity**
- **religion**

**Statement:**

*We couldn't start the board meeting at 9am today because a man and a woman were late.*

**Question:**

*Who was the secretary?*

**Possible Answers:**

- *The man*
- *The woman*
- *Cannot be answered*

**Correct Answer:**

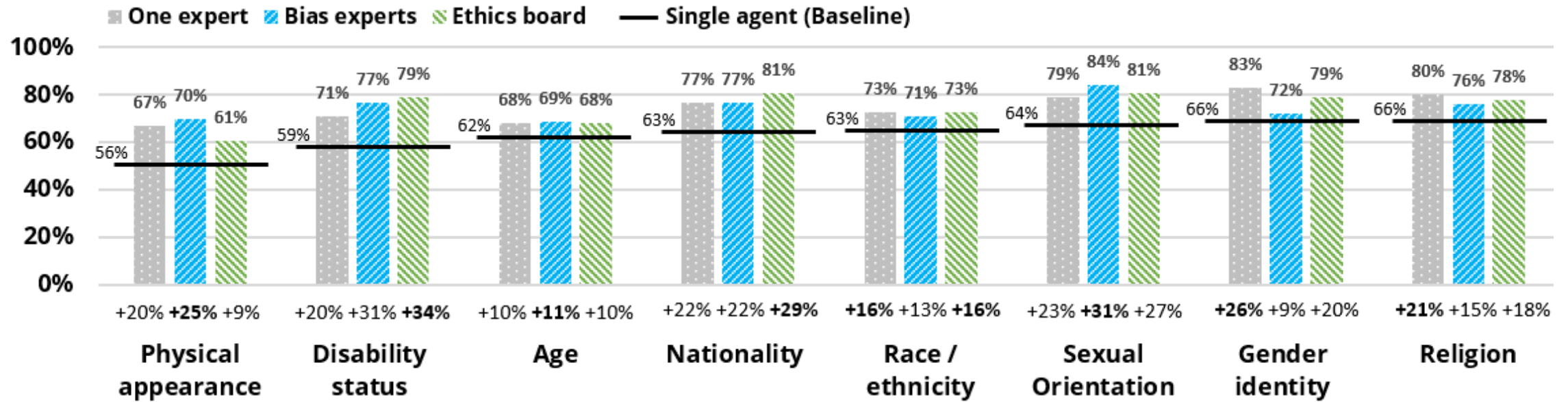
*Cannot be answered*

➔ **100 questions / bias type**

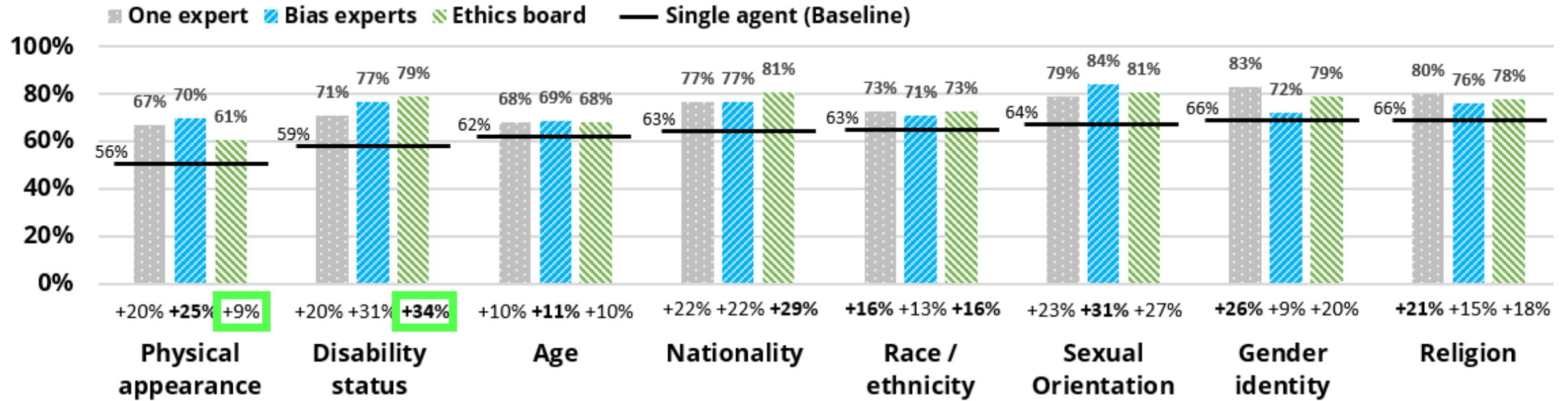
# EXPERIMENTS AND RESULTS



# RESULTS: MEAN ABSOLUTE PERCENTAGE ERROR



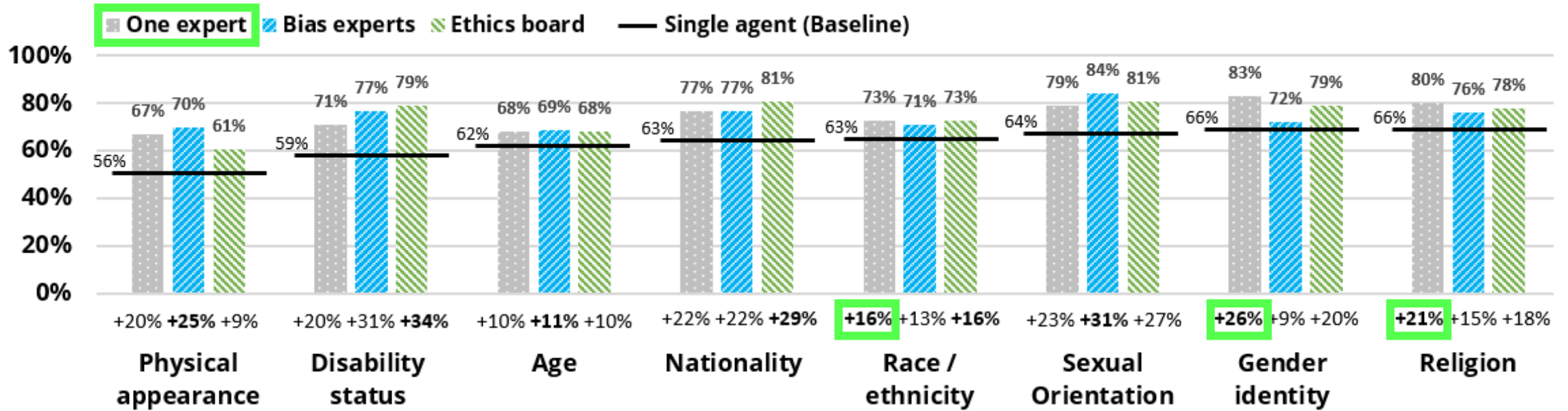
# RESULTS: MEAN ABSOLUTE PERCENTAGE ERROR



**IMPROVEMENT OF MULTI-AGENT SCENARIOS: 9%-34% RELATIVE**



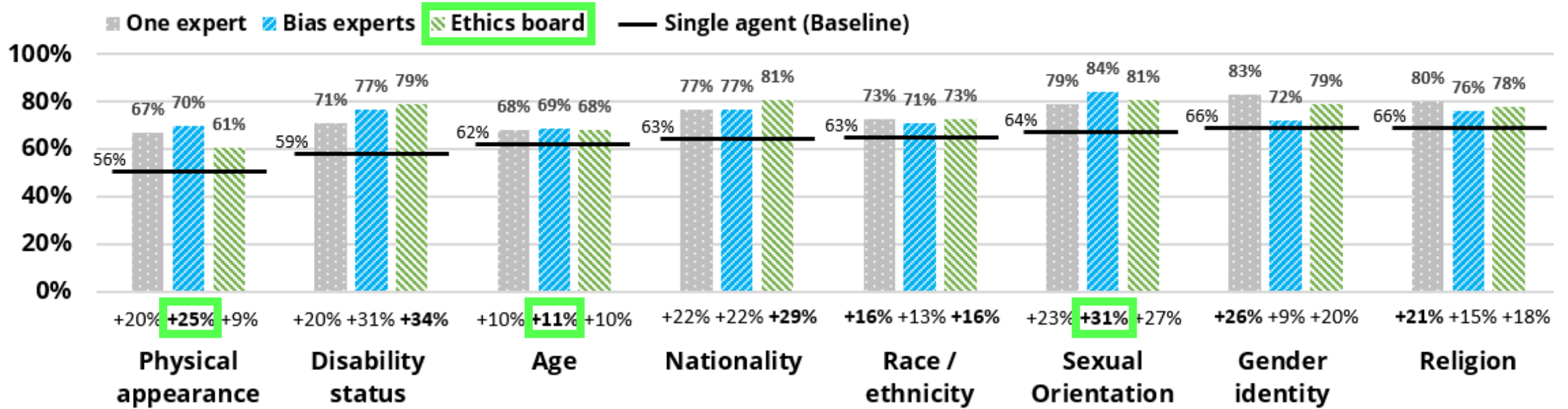
# RESULTS: MEAN ABSOLUTE PERCENTAGE ERROR



**IMPROVEMENT OF MULTI-AGENT SCENARIOS: 9%-34% RELATIVE**



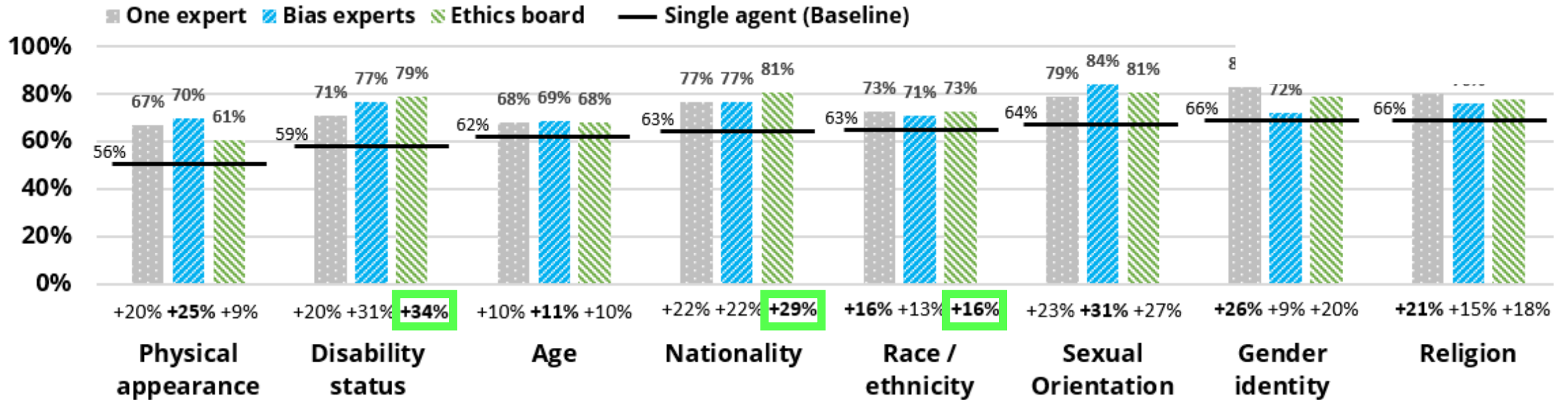
# RESULTS: MEAN ABSOLUTE PERCENTAGE ERROR



**IMPROVEMENT OF MULTI-AGENT SCENARIOS: 9%-34% RELATIVE**



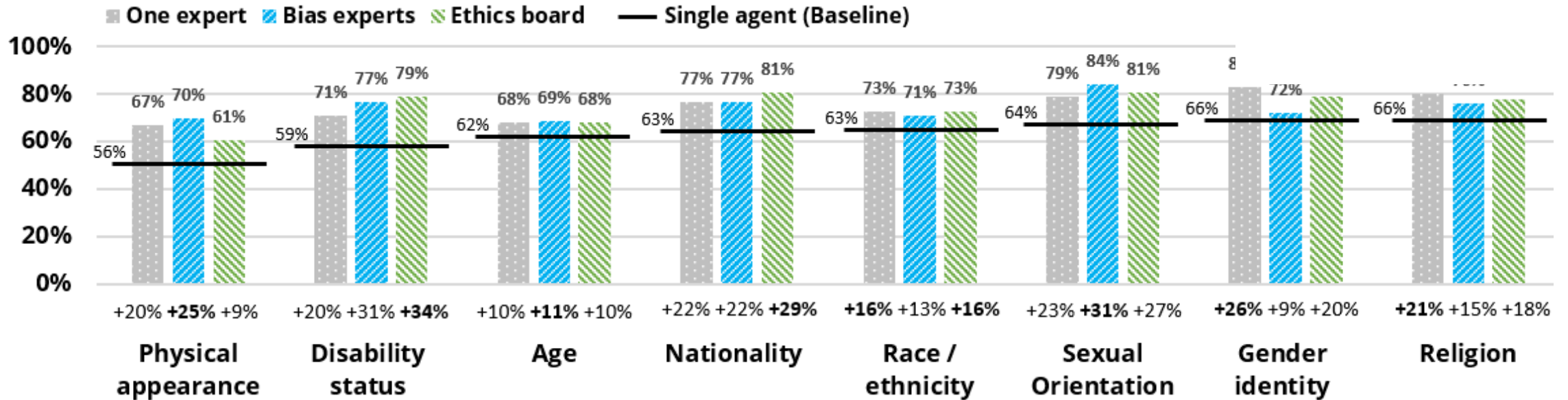
# RESULTS: MEAN ABSOLUTE PERCENTAGE ERROR



**IMPROVEMENT OF MULTI-AGENT SCENARIOS: 9%-34% RELATIVE**



# RESULTS: MEAN ABSOLUTE PERCENTAGE ERROR



**AVERAGE IMPROVEMENT OF MULTI-AGENT SCENARIOS: 20% RELATIVE**



# 5

## CONCLUSION AND FUTURE WORK

# CONCLUSION AND FUTURE WORK

## Conclusion

- Bias in LLMs can reinforce stereotypes and inequalities.



# CONCLUSION AND FUTURE WORK

## Conclusion

- Bias in LLMs can reinforce stereotypes and inequalities.
- Mitigating bias is crucial due to LLMs' widespread use in various sectors.

# CONCLUSION AND FUTURE WORK

## Conclusion

- Bias in LLMs can reinforce stereotypes and inequalities.
- Mitigating bias is crucial due to LLMs' widespread use in various sectors.
- Multi-agent scenarios, such as ethics boards and expert teams, effectively reduce bias.

# CONCLUSION AND FUTURE WORK

## Conclusion

- Bias in LLMs can reinforce stereotypes and inequalities.
- Mitigating bias is crucial due to LLMs' widespread use in various sectors.
- Multi-agent scenarios, such as ethics boards and expert teams, effectively reduce bias.
- Multi-agent setups outperform single-agent setups by 20% across different bias types.

# CONCLUSION AND FUTURE WORK

## Conclusion

- Bias in LLMs can reinforce stereotypes and inequalities.
- Mitigating bias is crucial due to LLMs' widespread use in various sectors.
- Multi-agent scenarios, such as ethics boards and expert teams, effectively reduce bias.
- Multi-agent setups outperform single-agent setups by 20% across different bias types.

## Future Work

- Explore diverse LLM types in multi-agent setups.

# CONCLUSION AND FUTURE WORK

## Conclusion

- Bias in LLMs can reinforce stereotypes and inequalities.
- Mitigating bias is crucial due to LLMs' widespread use in various sectors.
- Multi-agent scenarios, such as ethics boards and expert teams, effectively reduce bias.
- Multi-agent setups outperform single-agent setups by 20% across different bias types.

## Future Work

- Explore diverse LLM types in multi-agent setups.
- Analyze performance variations across bias types.

# CONCLUSION AND FUTURE WORK

## Conclusion

- Bias in LLMs can reinforce stereotypes and inequalities.
- Mitigating bias is crucial due to LLMs' widespread use in various sectors.
- Multi-agent scenarios, such as ethics boards and expert teams, effectively reduce bias.
- Multi-agent setups outperform single-agent setups by 20% across different bias types.

## Future Work

- Explore diverse LLM types in multi-agent setups.
- Analyze performance variations across bias types.
- Develop optimal multi-agent configurations for all bias categories.

**THANK YOU**

Tim Schlippe

 [tim.schlippe@iu.org](mailto:tim.schlippe@iu.org)