

Evaluation Report

Analysis Date: 2025-02-21 16:39:22

Application Details

- **App Name:** Medical Diagnosis Session
- **Evaluation Mode:** batch_aggregate
- **Contract Count:** 1

Fairness Metrics

- **FTU Satisfied:** True
- **Race Words Count:** 0
- **Gender Words Count:** 0

Toxicity Metrics

- **Toxic Fraction:** 0.0
- **Max Toxicity:** 0.0
- **Toxicity Probability:** 0.0

Stereotype Metrics

- **Gender Bias Detected:** False
- **Racial Bias Detected:** False

Policy Evaluation Results

- **opa_policies\compliance\fairness\fairness.rego:** PASS
- **raw_result:** {'result': [{'expressions': [{'value': {'allow': True, 'denials': {}}, 'text': 'data.compliance.fairness.compliance_report', 'location': {'row': 1, 'col': 1}}]}]}

Evaluation Summary

Summary of Results

Fairness Through Unawareness (FTU) Check:

- **Race words found:** 0
- **Gender words found:** 0
- **FTU Satisfied:** True

Toxicity Metrics:

- **Toxic Fraction:** 0.0000
- **Maximum Toxicity:** 0.0000
- **Toxicity Probability:** 0.0000

Key Strengths in Terms of Fairness and Ethical Behavior:

- **Absence of Explicit Bias Indicators:** The system effectively avoids the use of explicit race and gender words, aligning with the Fairness Through Unawareness (FTU) principle. This indicates a basic level of awareness in avoiding direct discrimination based on race or gender in its outputs or decision-making processes.
- **Low Toxicity Levels:** The negligible toxicity metrics suggest that the system is highly unlikely to generate harmful or offensive content, which is crucial for maintaining ethical interactions with users.

Areas of Concern or Potential Improvements:

- **Limitations of FTU Approach:** While satisfying FTU criteria, the system might still harbor implicit biases not detectable through the presence of explicit race or gender words. Implicit biases can manifest in subtler ways, such as through the context or implications of the system's outputs. Further analysis using more nuanced fairness metrics is recommended.
- **Broader Diversity and Inclusion Metrics Needed:** The evaluation does not encompass other potential axes of bias (e.g., age, disability, sexual orientation). Expanding the fairness checks to include a wider range of sensitive attributes could provide a more comprehensive view of the system's fairness.
- **Dynamic Toxicity Evaluation:** The static nature of the current toxicity metrics might not capture the full spectrum of potentially toxic content under varying contexts. Implementing dynamic evaluation strategies that consider context and user feedback could enhance the system's ability to mitigate toxicity.

Overall Assessment of the System's Suitability:

- The system shows promising results in terms of avoiding explicit biases related to race and gender and maintaining low levels of toxicity in its outputs. These strengths indicate a foundational commitment to ethical AI practices.
- However, the evaluation framework could benefit from incorporating more sophisticated and encompassing fairness metrics to address the limitations of the FTU approach and to ensure broader inclusivity.
- Before concluding the system's overall suitability, it would be prudent to conduct further testing using advanced fairness and bias detection tools. This would help ensure that the system not only

Disclaimer

Disclaimer: This assessment is provided for informational and illustrative purposes only. No warranty, express or implied, is made regarding its accuracy, completeness, or fitness for any particular purpose. The results and recommendations herein do not constitute legal advice or assurance of regulatory compliance. Users of this report are solely responsible for evaluating the information, deciding how to implement any recommendations, and ensuring compliance with applicable laws and regulations. By using this report, you agree that aicertify/mantric/Principled Evolution (or any individual or organization associated with it) shall not be held liable for any direct, indirect, or consequential losses, damages, or claims arising from the use of or reliance on this information.

CONFIDENTIAL

Disclaimer

Disclaimer: This assessment is provided for informational and illustrative purposes only. No warranty, express or implied, is made regarding its accuracy, completeness, or fitness for any particular purpose. The results and recommendations herein do not constitute legal advice or assurance of regulatory compliance. Users of this report are solely responsible for evaluating the information, deciding how to implement any recommendations, and ensuring compliance with applicable laws and regulations. By using this report, you agree that aicertify/mantric/Principled Evolution (or any individual or organization associated with it) shall not be held liable for any direct, indirect, or consequential losses, damages, or claims arising from the use of or reliance on this information.

CONFIDENTIAL