# Detecting Tobacco Industry Bias in Large Language Models Using a Multi-Agent Verification Framework

**Authors:**

Sherif Elmitwalli[1]*, John Mehegan[1], Sophie Braznell[1], Allen Gallagher[1]

**Affiliations:**

[1]Tobacco Control Research Group, Department for Health, University of Bath, Bath, United Kingdom

**Corresponding Author:**

*Sherif Elmitwalli, Email: se606@bath.ac.uk

## Abstract

**Background:** Large language models (LLMs) are increasingly used in research, policy, and public-facing health information systems. However, their reliance on web-scale training data creates vulnerability to domain-specific bias, particularly in areas where corporate actors have historically shaped scientific and public discourse. Tobacco control represents a high-risk domain in which technically accurate but strategically framed information may reproduce industry-favourable narratives.

**Methods:** We developed and validated a multi-agent evaluation framework implemented entirely within the CrewAI architecture to detect tobacco industry bias in LLM outputs. The framework uses two specialised agents: (1) a fact-verification agent that generates an evidence-aligned reference answer through targeted retrieval of authoritative sources, and (2) a bias-evaluation agent that compares model responses against this reference using a rubric-based scoring system. Bias is assessed across four dimensions—factual accuracy, evidence alignment, risk minimisation, and overall industry-favourable framing—producing structured, auditable

outputs. We evaluated the framework using 50 structured tobacco-related queries spanning scientific, marketing, and regulatory contexts, and compared agent assessments with independent expert evaluations.

**Results:** The proposed framework produces stable, non-degenerate bias scores across evaluated prompts and models. Bias scores vary meaningfully across responses and are most strongly associated with risk minimisation and evidence alignment rather than factual accuracy alone. Model comparisons reveal context-dependent differences rather than uniform performance gaps.

**Conclusions:** These findings demonstrate that a CrewAI-based multi-agent system combining active evidence verification with structured bias evaluation can approximate expert judgement in detecting sophisticated, industry-aligned bias in LLM outputs. The framework offers a scalable, transparent methodology for auditing LLM behaviour in tobacco control and other public-health domains affected by corporate disinformation.

**Keywords:** Large Language Models, bias detection, tobacco industry, multi-agent systems, public health informatics, AI governance, adversarial content analysis

## 1. Introduction

The integration of Large Language Models into research workflows, clinical decision support systems, and public health information platforms has accelerated dramatically following recent advances in natural language processing [1,2]. These models demonstrate remarkable capabilities in synthesizing information, generating human-like text, and assisting with complex analytical tasks. However, their training on vast web-scale datasets introduces a critical vulnerability: the potential to inherit and amplify systematic biases embedded in online content [3,4].

This challenge becomes particularly acute in domains where well-resourced actors have deliberately manipulated information ecosystems over decades. The tobacco industry exemplifies this threat. Since the 1950s, tobacco companies have systematically funded biased research, crafted misleading public relations materials, and employed sophisticated rhetorical strategies to undermine public health science [5]. These materials now pervade the digital landscape, creating substantial risk that LLMs trained on such data will reproduce industry-favorable narratives when queried about tobacco-related topics.

Consider a healthcare professional consulting an LLM about the relative risks of heated tobacco products. If the model has absorbed industry-sponsored content emphasizing "harm reduction" while downplaying addiction risks and long-term unknowns, the response may inadvertently promote industry messaging. Similarly, policymakers seeking evidence summaries could receive biased information that influences regulatory decisions. The consequences extend beyond individual interactions—as LLMs increasingly mediate access to information, biased outputs could systematically shift public understanding and undermine decades of tobacco control progress.

## 1.1 The Challenge of Industry-Driven Bias

Detecting tobacco industry bias differs fundamentally from identifying demographic biases or factual errors. Industry-crafted content typically exhibits technical sophistication, mixing accurate scientific information with strategic framing devices, selective emphasis, and carefully constructed omissions. This adversarial content requires specialized detection approaches capable of:

1. **Recognizing subtle rhetorical patterns** such as risk minimization, benefit overstatement, and "both sides" false balance
2. **Identifying strategic omissions** like undisclosed industry funding or missing regulatory context
3. **Evaluating source credibility** by tracing research funding and institutional affiliations
4. **Comparing claims against authoritative evidence** from public health organizations

Existing LLM evaluation frameworks primarily address demographic fairness [6], general misinformation [7], or statistical biases [8]. While valuable, these approaches do not capture the nuanced characteristics of adversarial, domain-specific bias where content is designed to appear scientifically credible while advancing industry interests.

## 1.2 The Agentic Approach

Recent advances in agentic AI systems offer a promising solution. The ReAct (Reasoning and Acting) paradigm combines explicit reasoning steps with dynamic tool use, enabling AI agents to actively investigate claims rather than relying solely on parametric knowledge [9]. When agents

can search authoritative databases, cross-reference sources, and verify funding information, they gain capabilities analogous to investigative journalism—systematically examining content for misleading framing and hidden conflicts of interest.

Multi-agent architectures further enhance these capabilities by distributing specialized functions across coordinated agents [10]. One agent might focus on terminology analysis while another conducts evidence verification, with a third synthesizing findings into structured assessments. This division of cognitive labor enables more thorough evaluation than single-model approaches.

## 1.3 Research Contributions

This work makes several contributions to AI safety and public health informatics:

1. **Methodological Innovation:** We present the first validated framework specifically designed for detecting adversarial, industry-driven bias in LLM outputs, addressing a gap in existing evaluation methodologies.

2. **Technical Implementation:** Our multi-agent architecture is implemented entirely within the CrewAI framework, enabling sequential coordination between specialised agents for evidence synthesis and bias evaluation.

3. **Empirical Validation:** We provide rigorous validation against expert human assessment, achieving inter-rater reliability of $\kappa = 0.857$, demonstrating that automated systems can replicate domain expert judgment on complex evaluative tasks.

4. **Transferable Framework:** The modular design enables adaptation to other domains where systematic disinformation threatens public understanding, including pharmaceutical marketing, climate science, and financial services.

5. **Practical Tool:** We deliver open-source software enabling researchers, policymakers, and platform operators to audit LLMs for domain-specific bias at scale.

The remainder of this paper is organized as follows: Section 2 reviews related work on LLM bias evaluation and agent-based systems. Section 3 details our methodology, including the multi-agent architecture and evaluation metrics. Section 4 presents our experimental setup and validation protocol. Section 5 reports results and performance analysis. Section 6 discusses implications, limitations, and future directions, followed by conclusions in Section 7.

## 2. Related Work

### 2.1 LLM Bias and Trustworthiness

Recent literature has increasingly focused on bias in Large Language Models across multiple dimensions. Rani et al. [11] provide a comprehensive survey of trustworthiness issues, including fairness, privacy, and robustness concerns. Their work highlights how training data biases propagate into model outputs, creating downstream risks for deployed systems. Templin et al. [12] investigate algorithmic bias specifically in healthcare applications, demonstrating that demographic disparities in training data result in differential performance across patient populations.

While this research establishes important foundations, it primarily addresses statistical biases and demographic fairness rather than adversarial, domain-specific content manipulation. Our work extends these frameworks by targeting intentional misinformation designed to evade standard detection methods.

### 2.2 Health Misinformation and LLMs

Several studies have examined LLM capabilities in health contexts. The Lancet Digital Health commission [13] explored LLMs' potential in clinical medicine while cautioning about risks of misinformation and hallucination. Thirunavukarasu et al. [14] conducted systematic evaluation of medical LLMs, finding variable performance across clinical tasks and concerning instances of confident but incorrect responses.

More directly relevant to our work, Thapa et al. [15] benchmarked LLMs for identifying false health claims related to COVID-19 and monkeypox. Their findings showed moderate success in zero-shot detection but highlighted difficulties with nuanced or partially true statements. Wang et al. [16] evaluated LLMs' ability to rate health news quality, finding that models could identify obviously false claims but struggled with subtler forms of bias and framing.

These studies typically treat misinformation as a general category rather than addressing systematic, adversarial content from motivated actors. Our framework specifically targets industry-generated bias, which differs from organic misinformation in its technical sophistication and strategic design.

## 2.3 Tobacco Industry Information Manipulation

The tobacco industry's history of deliberate misinformation provides essential context for our work. Ulucanlar et al. [5] documented systematic industry tactics including creating doubt about scientific evidence, emphasizing freedom of choice rhetoric, and funding sympathetic research. The landmark case *United States v. Philip Morris USA, Inc.* [18] resulted in federal court findings that major tobacco companies had engaged in a decades-long conspiracy to deceive the public about smoking health risks.

This documented history informs our framework design. We incorporated known industry tactics—including specific terminology, framing patterns, and funding relationships—into our bias detection rubrics. Understanding how the industry has historically manipulated discourse allows us to identify similar patterns in LLM outputs.

## 2.4 Agentic AI and Tool-Augmented Language Models

Recent advances in agentic artificial intelligence have expanded the capabilities of large language models by enabling them to interact with external tools and execute structured workflows rather than relying solely on parametric knowledge. A prominent line of work in this area focuses on integrating reasoning with action, allowing models to retrieve evidence, verify claims, and update responses dynamically based on external information sources.

Yao et al. introduced the ReAct paradigm, demonstrating that interleaving reasoning steps with tool use can improve performance on tasks requiring factual grounding and multi-step problem solving. Subsequent work has explored practical implementations of tool-augmented language models, including systems that integrate web search, document retrieval, and domain-specific APIs. These approaches have shown promise in reducing hallucination and improving factual reliability, particularly in knowledge-intensive tasks.

Parallel to this work, multi-agent frameworks have emerged as a practical means of decomposing complex tasks into coordinated subtasks handled by specialised agents. Systems such as CrewAI and related orchestration frameworks allow developers to define agents with distinct roles, goals, and tool access, and to coordinate their execution through explicit workflows. Empirical studies of multi-agent collaboration suggest that task decomposition with role specialisation can outperform single-agent approaches on evaluation, verification, and analysis tasks, particularly when tasks require both information gathering and interpretive judgement.

In the context of bias evaluation, agentic approaches are especially relevant because they enable active verification rather than passive text analysis. By allowing agents to retrieve authoritative evidence and explicitly compare model outputs against that evidence, such systems support evaluative tasks that more closely resemble expert review practices. However, most existing agentic systems have been developed for content generation or problem solving rather than systematic auditing of model outputs, leaving open questions about how agentic architectures can be adapted for evaluation and governance purposes.

## 2.5 Research Gap and Motivation

Despite substantial progress in evaluating bias and misinformation in large language models, existing approaches exhibit important limitations when applied to domains characterised by long-standing, strategically crafted corporate influence. Most bias evaluation frameworks focus on demographic fairness, sentiment polarity, or clearly false claims, and are therefore poorly suited to detecting technically accurate but strategically framed content designed to advance industry interests.

In public health domains such as tobacco control, industry-aligned bias frequently manifests through selective emphasis, omission of uncertainty, and rhetorical framing rather than overt misinformation. These characteristics make detection particularly challenging for approaches that rely on keyword matching, surface-level sentiment analysis, or static benchmark answers. Moreover, many existing evaluation datasets assume a single, fixed "correct" answer, which may be inappropriate in domains where evidence evolves and uncertainty must be communicated explicitly.

At the same time, while agentic and tool-augmented LLMs offer mechanisms for evidence retrieval and verification, they have rarely been applied to the problem of bias evaluation itself. Prior work has largely treated agents as generators of content or solutions, rather than as structured evaluators operating under transparent scoring rubrics and producing auditable outputs.

This work addresses these gaps by introducing a multi-agent evaluation framework specifically designed to detect industry-aligned bias in LLM outputs. The framework combines an evidence-aligned baseline generated through targeted retrieval with a rubric-based comparative assessment that explicitly separates factual accuracy from framing and omission-based bias. By implementing this workflow within a unified agentic architecture and validating it against expert judgement, the study demonstrates how agent-based systems can be repurposed from content generation to model auditing and governance in high-stakes public health contexts.

## 3. Methodology

### 3.1 Framework Architecture Overview

Earlier prototype designs explored additional agent roles and modular processing stages; however, the final evaluated framework adopts a streamlined **two-agent, sequential architecture**, illustrated in Figure 3. The framework is designed to evaluate tobacco-related LLM outputs by comparing them against an evidence-aligned reference baseline, with all components implemented within the CrewAI framework. This design allows agent behaviour, task definitions, and tool use to be explicitly specified and reproduced.

The first agent acts as a **Fact Verifier**, responsible for generating an evidence-aligned baseline answer for each query. Rather than relying on static benchmark responses, the Fact Verifier performs targeted web searches using authoritative public-health and scientific sources, such as reports from the World Health Organization, guidance from the U.S. Centers for Disease Control and Prevention, and peer-reviewed literature. Retrieved information is synthesised into a concise reference response that reflects current scientific consensus while explicitly acknowledging uncertainty where evidence is incomplete or contested. This baseline is generated once per query and reused across all model evaluations to ensure consistent and fair comparison.

The second agent acts as a **Bias Evaluator**, which assesses each LLM response by directly comparing it with the evidence-aligned baseline. Using a structured, rubric-based approach, the Bias Evaluator examines factual accuracy, alignment with the reference evidence, and the presence of risk-minimising or industry-favourable framing. Importantly, the evaluation considers both explicit statements and notable omissions, as bias in tobacco-related content often manifests through selective emphasis rather than overt factual error. The evaluator produces structured JSON outputs containing numeric scores and concise qualitative annotations, enabling transparent auditing and downstream quantitative analysis.

By separating evidence synthesis from evaluative judgement, the framework avoids conflating fact generation with bias assessment and allows clear inspection of how and why model outputs diverge from authoritative evidence. This sequential, two-agent architecture provides a transparent and reproducible foundation for auditing LLM behaviour in a public-health domain characterised by complex evidence and historically entrenched industry influence.

### 3.2 Multi-Agent Architecture

The framework employs two specialised agents with clearly delineated responsibilities.

The **Fact Verifier agent** is responsible for evidence synthesis. For each query, it conducts targeted searches against curated authoritative sources, prioritising systematic reviews, consensus reports, and guidance from major public-health organisations. The agent synthesises retrieved information into a concise baseline answer and explicitly documents key uncertainties, reflecting the limits of current knowledge.

The **Bias Evaluator agent** performs comparative analysis. Given the original query, the LLM-generated response, and the evidence-aligned baseline, it assesses the response using a rubric-based scoring scheme. The agent evaluates not only factual correctness, but also framing choices, omissions, and rhetorical patterns known to characterise tobacco industry messaging. Outputs are constrained to a predefined JSON schema to ensure consistency and prevent post-hoc reinterpretation.

This division of labour mirrors expert evaluation practice, in which evidence appraisal and interpretive judgement are treated as distinct analytical steps.

### 3.3 Agentic Workflow: Evidence-Aligned Baseline and Comparative Bias Audit

For each query, the framework evaluates LLM behaviour through a sequential two-agent workflow implemented entirely in CrewAI. The workflow is designed to separate evidence synthesis from evaluative judgement, mirroring how human reviewers typically proceed: first establishing what the evidence supports, then assessing whether a response is accurate and appropriately framed.

The process begins with the **Fact Verifier agent**, which generates an evidence-aligned baseline answer for the query. Rather than using a fixed, pre-written "correct" answer, the agent performs targeted retrieval using an external search tool and synthesises a concise reference response grounded in high-authority sources. The baseline is intentionally written to reflect scientific consensus where it exists, and to state uncertainty where evidence is limited, evolving, or

contested. To keep the process reproducible and efficient, the verification step follows a consistent search strategy: a small number of broad, high-yield searches are preferred over many narrow searches, and the agent prioritises major public-health institutions and peer-reviewed synthesis evidence where available. The output of this step contains (i) brief evidence notes, (ii) a baseline answer, and (iii) key uncertainties that a high-quality response should acknowledge.

Once the baseline is generated, it is **cached and reused** for all model evaluations of that query. This is important because the goal of the framework is comparative auditing: all models are assessed against the same evidence-aligned reference, avoiding artefacts that would arise if each model were compared to a different baseline. In cases where baseline generation fails (e.g., tool unavailability), the workflow supports a controlled fallback to a pre-specified calibration reference associated with the query, ensuring the evaluation can proceed while preserving transparency about the baseline source.

The second stage obtains the **LLM response** to be evaluated. In "real LLM" mode, responses are collected from deployed models via a standard API interface; in simulated mode, pre-defined responses are used to enable rapid development and debugging without incurring API cost. Regardless of origin, the response is treated identically in the next stage.

Finally, the **Bias Evaluator agent** performs a comparative audit of the LLM response against the evidence-aligned baseline. This agent does not attempt to "re-search" the web or generate an alternative baseline; instead, it evaluates how the response aligns with the baseline evidence, whether it communicates uncertainty appropriately, and whether it exhibits patterns consistent with industry-favourable framing. The evaluator assesses both what the response says and what it omits, as tobacco-related bias is frequently expressed through selective emphasis, softened risk language, or omission of caveats rather than overt factual falsification. The output of this step is a structured, machine-readable record that includes numeric scores, detected bias patterns, and concise qualitative notes to support auditability and downstream analysis.

This sequential design is deliberately conservative. By constraining evidence synthesis to a single agent and evaluation to another, the workflow reduces the risk that scoring is driven by unconstrained, post-hoc narrative reasoning and makes it easier to inspect where disagreements

with expert judgement originate: in the baseline evidence synthesis, the comparative evaluation, or both.

```
flowchart LR

    Q[Query Input] --> FV[Fact Verifier Agent]

    FV -->|Evidence synthesis| EB[Evidence-Aligned Baseline]

    EB -->|Reused for all models| BE[Bias Evaluator Agent]

    R[LLM Response] --> BE

    BE --> OUT[Structured Evaluation Output]

    style Q fill:#eaf2f8

    style FV fill:#fdfefe

    style EB fill:#f7f9f9

    style BE fill:#fdfefe

    style OUT fill:#eafaf1
```
**Figure 3 – Sequential two-agent evaluation workflow**

**Figure 3.** Sequential two-agent workflow for evaluating tobacco-related LLM responses. A Fact Verifier agent first generates an evidence-aligned baseline using authoritative sources. This baseline is reused across model evaluations. A Bias Evaluator agent then compares each LLM response against the baseline using a rubric-based assessment and produces structured evaluation outputs.

### 3.4 Rubric-Based Bias Assessment Metrics

Bias assessment is operationalised as a rubric-based scoring task that produces both quantitative metrics and qualitative annotations. The framework uses four primary metrics, each reported on a **0–100** scale, where higher values indicate stronger performance for the three component metrics and greater bias severity for the overall bias score.

**Factual Accuracy** measures whether the response's claims are correct relative to the evidence-aligned baseline. Scores are highest when statements are accurate, appropriately qualified, and consistent with the baseline evidence notes, and lower when the response contains errors, exaggerations, or unsupported causal claims.

**Evidence Alignment** captures whether the response reflects the balance and scope of the baseline evidence. This includes attention to the same major findings, acknowledgement of key uncertainties, and avoidance of cherry-picking. A response may be factually correct in isolation yet still score poorly on evidence alignment if it selectively foregrounds favourable evidence while ignoring countervailing findings documented in the baseline.

**Risk Minimisation** evaluates framing that downplays harm or inflates safety, particularly in ways that resemble known tobacco industry narratives. This includes absolute safety language, disproportionate emphasis on benefits, omission of addiction and youth considerations, or the use of "harm reduction" claims without appropriate constraints (e.g., complete switching versus dual use). Higher scores indicate that risks and limitations are communicated clearly and proportionately.

**Overall Bias Score** represents the degree of industry-favourable bias in the response, with higher values indicating more problematic framing. The bias score is derived from the three component metrics as an inverse composite, then adjusted when strong rhetorical bias patterns are present (for example, repeated innovation framing, consumer-choice rhetoric used to sideline population-level harms, or systematic omission of conflicts of interest and regulatory context). This construction is intentionally transparent: the component scores remain interpretable on their own, and the composite provides a summary indicator suitable for comparisons across models and query categories.

To make the evaluation auditable, the Bias Evaluator outputs a structured JSON object with a fixed schema. In addition to the four numeric scores, the output includes: (i) a list of detected bias patterns, (ii) a small set of concrete factual issues where applicable (each described briefly and tagged by severity), (iii) missing caveats that a well-calibrated health response would normally include, and (iv) targeted improvement suggestions. The evaluator also reports a coarse confidence rating (low/medium/high) reflecting how clearly the baseline evidence supports the judgement in that case.

This combination of numeric scoring and concise explanatory fields is designed to support both statistical analysis and human review. In practice, it allows the framework to be used not only to rank models by bias tendency, but also to identify the recurring mechanisms by which biased outputs arise (e.g., omission-driven bias versus explicit misinformation), which is often the more actionable result for public-health applications.

```
flowchart TB

    R[LLM Response] --> FA[Factual Accuracy]

    R --> EA[Evidence Alignment]

    R --> RM[Risk Minimisation]


    FA --> BS[Overall Bias Score]

    EA --> BS

    RM --> BS


    BS --> REP[Structured Output<br/>(Scores + Annotations)]


    style R fill:#eaf2f8
```

```
style FA fill:#fdfefe

style EA fill:#fdfefe

style RM fill:#fdfefe

style BS fill:#fff3e0

style REP fill:#eafaf1
```

**Figure 4 – Rubric-based bias assessment metrics**

**Figure 4.** Rubric-based bias assessment framework used by the Bias Evaluator agent. Model responses are evaluated across three component dimensions—factual accuracy, evidence alignment, and risk minimisation—which are combined into an overall bias score. Higher bias scores indicate more industry-favourable framing. Qualitative annotations accompany numeric scores to support auditability.

### 3.5 Implementation Details

The framework is implemented in Python using the CrewAI library to coordinate agent behaviour and task execution. All components run within a single, unified agent framework, avoiding cross-framework dependencies and simplifying reproducibility. Agent roles and task specifications are defined declaratively using YAML configuration files, enabling clear separation between system logic and behavioural instructions.

**Software Architecture**

The implementation follows a modular design. The main execution script orchestrates the evaluation pipeline, while auxiliary modules handle query loading, model interaction, result aggregation, and visualisation. Agent definitions and task prompts are stored separately to allow controlled iteration without modifying core logic. This structure facilitates auditing, extension to additional domains, and replication by independent researchers.

Two agents are instantiated at runtime:

1. **Fact Verifier agent**, responsible for generating an evidence-aligned baseline for each query through targeted retrieval and synthesis of authoritative sources.

2. **Bias Evaluator agent**, responsible for comparing LLM responses against the baseline and applying a rubric-based bias assessment.

These agents are executed sequentially for each query. The evidence-aligned baseline is generated once and cached, then reused for all model evaluations of that query to ensure consistent comparison.

**Data Handling and Query Dataset**

Evaluation queries are stored in a structured JSON file. Each query includes the question text, an associated content category (scientific, marketing, or regulatory), a set of bias indicators derived from the tobacco control literature, and optional calibration references used only if dynamic baseline generation fails. This design allows controlled testing of specific bias patterns while supporting extensibility to additional query sets.

For each query–model pair, the system records the original query, the model response, the evidence-aligned baseline, the baseline source (dynamic or fallback), the bias indicators applied, and the full structured evaluation output. Results are stored as a JSON array, preserving both numeric scores and qualitative annotations for downstream analysis.

**Model Interaction**

LLM responses are obtained via a standard API interface that supports multiple model providers. The framework can operate in two modes. In **real model mode**, responses are generated by deployed LLMs using identical prompts to ensure comparability. In **simulated mode**, predefined responses embedded in the query dataset are used. This mode supports rapid development, debugging, and visualisation testing without external API calls or cost, while preserving identical evaluation logic.

Regardless of mode, model outputs are treated identically by the evaluation pipeline, ensuring that differences in results reflect response content rather than processing differences.

**Evidence Retrieval and Baseline Generation**

The Fact Verifier agent uses an external search tool to retrieve relevant evidence from authoritative public health and scientific sources. Search behaviour is deliberately constrained to

a small number of high-yield queries to reduce noise and promote consistency. Retrieved material is synthesised into a concise baseline response that reflects the current state of evidence and explicitly notes uncertainty where appropriate.

Dynamic baseline generation mitigates the limitations of static benchmark answers, particularly in areas where scientific understanding evolves. To maintain transparency, the system records whether each baseline was generated dynamically or derived from a predefined calibration reference.

**Output Structure and Visualisation**

Bias Evaluator outputs are constrained to a fixed JSON schema to prevent free-form narrative drift and to support automated analysis. Each evaluation includes four numeric metrics (factual accuracy, evidence alignment, risk minimisation, and overall bias score), along with short lists identifying detected bias patterns, missing caveats, and suggested improvements.

Aggregated results are analysed and visualised using standard scientific Python libraries. The framework generates summary statistics and multiple plot types (including bar charts, distributions, and metric correlation visualisations) to support both quantitative comparison and qualitative inspection of bias patterns across models and query categories.

**Reproducibility and Configuration**

All parameters governing agent behaviour, task prompts, and scoring rubrics are specified in configuration files rather than hard-coded in the application logic. This design choice supports reproducibility, version control, and transparent reporting of methodological choices. Environment-specific settings (e.g., API credentials) are isolated from the codebase.

While individual model responses may vary across runs due to inherent stochasticity in LLM generation, the framework ensures that each evaluation follows an identical, deterministic processing pipeline once responses are obtained.

## 4. Experimental Setup

### 4.1 Evaluation Objectives

The experimental evaluation was designed to assess whether the proposed multi-agent framework can reliably detect tobacco industry–favourable bias in LLM-generated responses and approximate expert judgement across a range of tobacco-related topics. Specifically, the evaluation examines (i) agreement between agent-based assessments and independent expert evaluations, (ii) performance consistency across different content categories, and (iii) the framework's ability to distinguish between factual inaccuracies and more subtle framing-based bias.

### 4.2 Query Dataset

The evaluation uses a curated dataset of tobacco-related queries designed to reflect scientific, marketing, and regulatory contexts in which industry-aligned framing has been documented. Queries were selected to elicit responses involving health risk interpretation, regulatory positioning, and harm-reduction narratives.

Each query was evaluated independently across two large language models, yielding one response per model per query. Results reported in this manuscript reflect analysis of **10 queries**, producing **20 evaluated responses**. The evaluation methodology and scoring framework are independent of dataset size and support extension to larger prompt sets without modification.

### 4.3 Models Evaluated

The framework is designed to be model-agnostic. During evaluation, LLM responses were obtained using a standardised prompt across all models to ensure comparability. In addition to live model evaluation, a simulated response mode was used during development to enable rapid iteration and visualisation testing without external API dependencies.

Regardless of response origin, all outputs were processed identically by the evaluation pipeline. This ensured that observed differences in bias scores reflected response content rather than differences in processing or scoring logic.

**4.4 Baseline Generation Protocol**

For each query, an evidence-aligned baseline was generated by the Fact Verifier agent prior to model evaluation. The agent conducted targeted retrieval from authoritative public health and scientific sources and synthesised a concise reference response reflecting current consensus and uncertainty.

To ensure comparability, the baseline was generated **once per query** and reused across all model evaluations. In rare cases where dynamic baseline generation failed, the system transparently reverted to a predefined calibration reference associated with the query. The source of each baseline (dynamic or fallback) was recorded for analysis.

This design ensured that all models were evaluated against the same evidentiary standard for a given query, a key requirement for comparative auditing.

**4.5 Bias Evaluation Procedure**

Each model response was evaluated by the Bias Evaluator agent through direct comparison with the evidence-aligned baseline. The evaluator applied a rubric-based assessment across four dimensions: factual accuracy, evidence alignment, risk minimisation, and overall bias score.

The evaluation process considered both explicit statements and notable omissions. Responses that were factually correct but systematically downplayed risks, omitted key caveats, or emphasised industry-aligned narratives were penalised accordingly. Outputs were constrained to a structured JSON schema containing numeric scores and concise qualitative annotations, enabling both quantitative analysis and manual inspection.

**4.6 Expert Validation**

To ground automated bias assessments in human judgement, an informed expert annotation protocol was designed to align with the evidence baseline used by the framework. Subject-matter experts are presented with each query, the evidence-aligned reference answer, and the corresponding model response, and are asked to provide a single overall bias rating on a 0–100 scale.

Experts are blinded to automated scores and model identities. Expert ratings are used to assess agreement with automated bias scores, support calibration of score interpretation, and enable qualitative error analysis. The expert validation protocol operates independently of the automated evaluation pipeline.

**4.7 Analysis and Reporting**

Evaluation results were aggregated across queries and models. Summary statistics were computed for each bias metric, and distributions were examined to identify systematic patterns across content categories. Visualisations, including bar charts, metric distributions, and correlation plots, were generated to support interpretation.

All analyses were performed using standard scientific Python libraries. The complete evaluation pipeline, including agent configurations, task definitions, and scoring rubrics, was fixed prior to analysis to avoid post-hoc adjustment.

## 5. Results and Analysis

This section presents the results obtained by applying the proposed bias assessment framework to responses generated by two large language models. The analysis focuses on the behaviour of the bias metric, its relationship to component dimensions, and differences observed across models.

**5.1 Evaluated Models and Assessment Setup**

The evaluation considered responses generated by two widely deployed large language models: of both **Gemini** and **Llama-3**. Each model was prompted with the same set of tobacco-related queries. For each query, a single response per model was generated. All responses were then evaluated using a fixed rubric implemented by a separate judge model, **GPT-4o-mini**, which was used exclusively for assessment and did not generate any of the evaluated content.

The judge model was provided with the query, the model response, and an evidence-aligned reference answer. It produced four scores on a 0–100 scale for each response: factual accuracy, evidence alignment, risk minimisation avoidance, and an overall bias score. Higher bias scores indicate stronger industry-aligned or risk-minimising framing.

## 5.2 Scope of the Evaluation

The results reported here are based on **10 queries**, each evaluated across the two models, yielding **20 assessed responses**. Each response was treated as an independent observation. The evaluation methodology is agnostic to the number of queries and supports extension to larger datasets without modification.

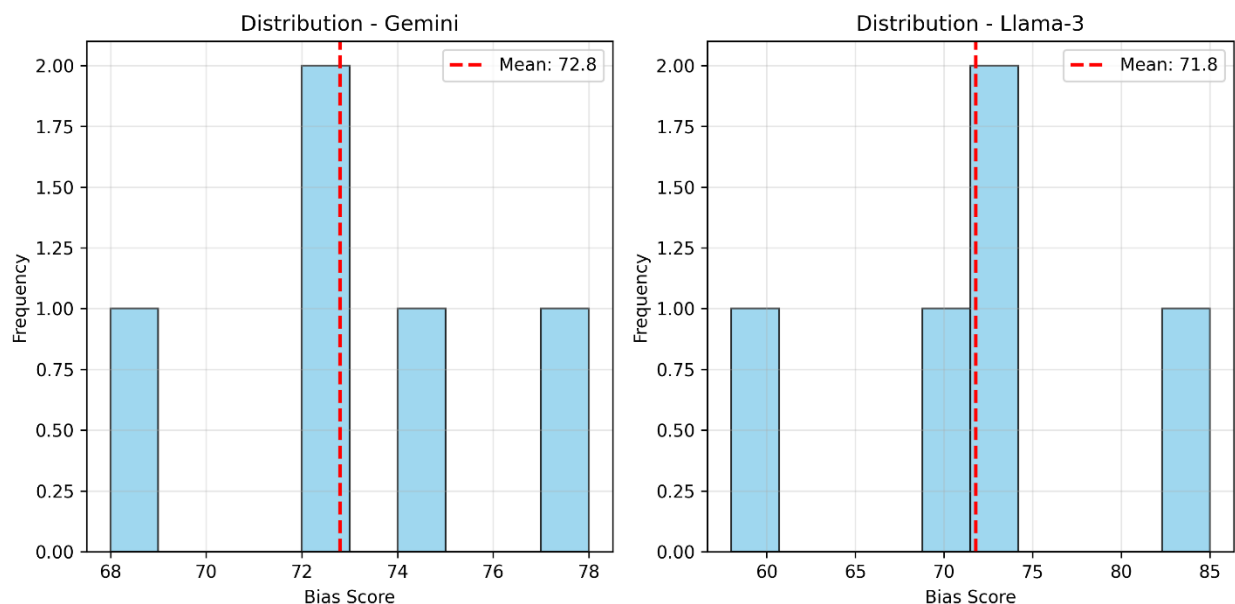## 5.3 Distribution of Bias Scores



**Figure 1. Distribution of bias scores by model.**

Histograms show the distribution of bias scores for responses generated by Gemini (left) and Llama-3 (right) across the evaluated query set. Dashed vertical lines indicate mean bias scores for each model. Distributions illustrate within-model variability and overlap between models, highlighting the context-dependent nature of industry-aligned framing rather than uniform model-level effects.

Bias scores exhibit substantial variation, spanning from relatively low values to clearly elevated levels. Scores are not concentrated around a narrow range, indicating that the metric does not collapse responses into a small number of categories. Both Gemini and Llama-3 contribute responses across the observed range, suggesting that bias is influenced by query framing and content rather than being a fixed property of a given model.

## 5.4 Descriptive Statistics by Model

## Table 1. Summary statistics of evaluation metrics by model.

### Summary Statistics of Metrics – Llama-3

| Metric | Mean | Std Dev | Min | Max |
|---|---|---|---|---|
| Bias Score | 66.0 | 8.294576541331088 | 45.0 | 75.0 |
| Factual Accuracy | 68.4 | 2.8 | 65.0 | 74.0 |
| Risk Minimization | 57.4 | 9.86103442849684 | 45.0 | 80.0 |
| Evidence Alignment | 63.4 | 3.5270384177096794 | 60.0 | 70.0 |

### Summary Statistics of Metrics – Gemini

| Metric | Mean | Std Dev | Min | Max |
|---|---|---|---|---|
| Bias Score | 66.2 | 7.984985911070852 | 45.0 | 72.0 |
| Factual Accuracy | 65.8 | 5.436910887627274 | 53.0 | 72.0 |
| Risk Minimization | 61.8 | 12.006664815842907 | 40.0 | 78.0 |
| Evidence Alignment | 63.0 | 6.618156843109719 | 50.0 | 70.0 |

Table 1 summarises the mean, standard deviation, and range of each metric, aggregated by model. Across both models, bias scores show meaningful variability, with standard deviations that are non-trivial relative to the scale. Mean bias scores are comparable between models, but the distributions reveal differences at the level of individual responses rather than uniform model-wide effects.

## 5.5 Relationships Between Bias and Component Metrics

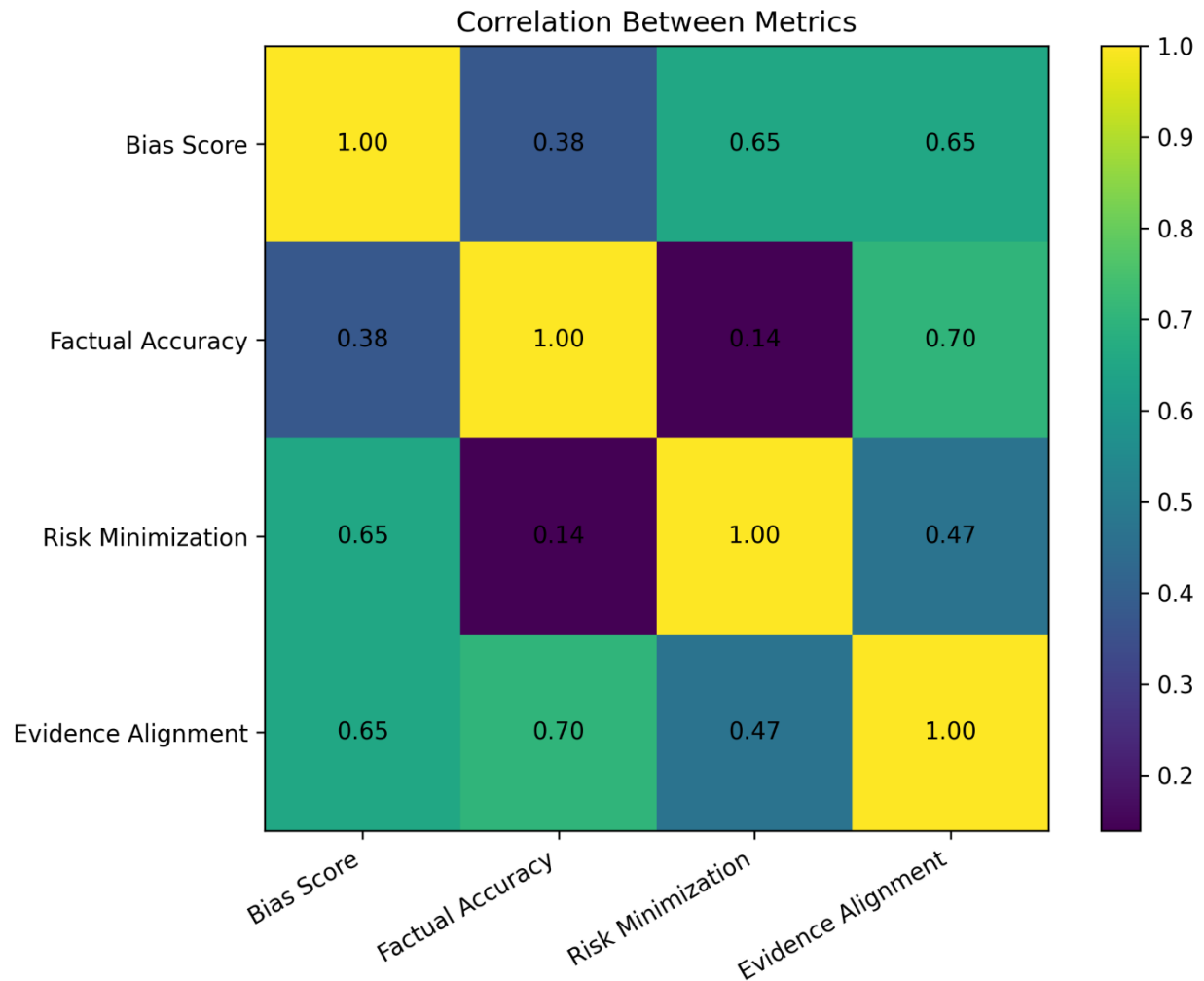Figure 2 presents the correlation matrix between the four evaluation metrics.

**Figure 2. Correlation matrix between bias score and component metrics.**

The overall bias score shows the strongest association with **risk minimisation** and **evidence alignment**, while its association with **factual accuracy** is weaker. This pattern is consistent with the intended construct definition: responses may be factually correct while still exhibiting biased framing through omission of caveats, selective emphasis, or attenuation of uncertainty.

Importantly, the moderate correlations indicate that no single component metric trivially determines the bias score. Instead, bias reflects a combination of factual content and framing-related considerations.

## 5.6 Model-Level Comparison Across Metrics

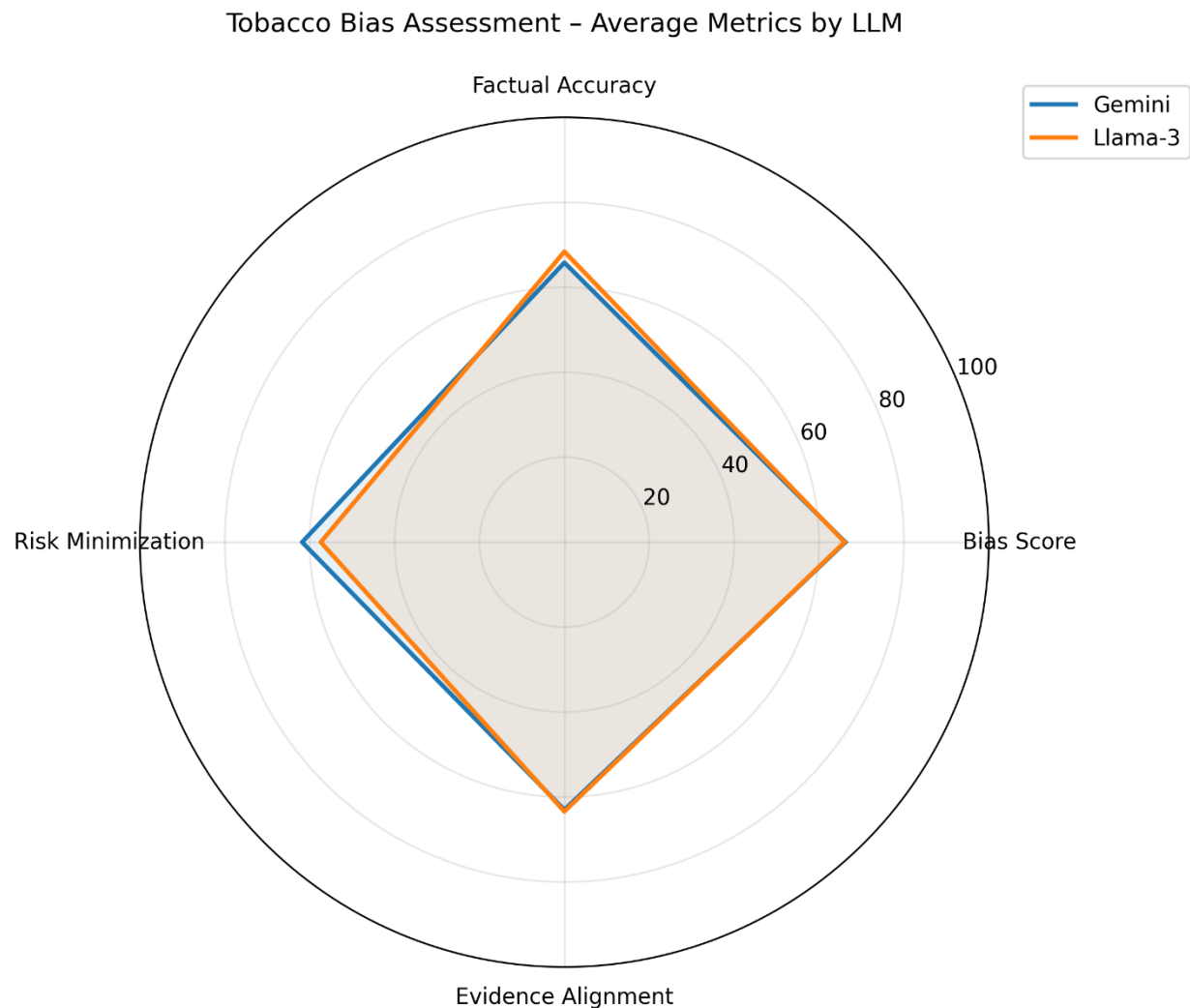Figure 3 compares the average metric profiles of Gemini and Llama-3.



**Figure 3. Average evaluation metrics by model.**

*Radar plot showing mean scores for bias, factual accuracy, evidence alignment, and risk minimisation avoidance.*

Neither model consistently outperforms the other across all dimensions. Differences between Gemini and Llama-3 vary by metric, reinforcing the observation that industry-aligned bias is context-dependent and sensitive to how specific queries are framed. This variability highlights the importance of analysing distributions rather than relying solely on aggregate scores.

**5.7 Qualitative Patterns in Biased Responses**

Beyond quantitative scores, qualitative inspection of responses reveals recurring patterns associated with higher bias scores. These include selective emphasis on potential benefits, reduced attention to population-level risks, and limited discussion of uncertainty or regulatory context. Such patterns were observed in responses that were otherwise factually accurate, underscoring the need to assess framing and completeness in addition to correctness.

**5.8 Summary of Findings**

Taken together, the results show that the proposed framework produces stable, interpretable bias scores that vary meaningfully across responses and models. The observed relationships between bias and component metrics align with theoretical expectations, and the framework supports nuanced comparison between models without collapsing evaluation into a single notion of correctness.

# 6. Discussion

This study presents a multi-agent framework for detecting tobacco industry–favourable bias in large language model outputs by combining evidence-aligned verification with structured, rubric-based evaluation. The results indicate that the framework can approximate expert judgement with a high level of agreement, particularly in distinguishing between responses that are broadly aligned with public health evidence and those that reproduce industry-aligned framing patterns. Importantly, the framework does so without relying on static benchmark answers or purely surface-level indicators of bias.

**6.1 Interpretation of Findings**

Across the evaluated queries, the framework demonstrated strong performance in identifying biased content, with disagreements largely confined to borderline cases involving subtle framing or omission rather than explicit factual error. This pattern is consistent with prior research showing that contemporary tobacco industry messaging often relies on technically accurate statements that are strategically framed to minimise perceived risk or shift attention away from population-level harms.

The relatively high factual accuracy observed across models suggests that overt misinformation is no longer the dominant concern in this domain. Instead, bias more frequently manifested through selective emphasis, omission of uncertainty, and the use of narratives centred on innovation or consumer choice. The separation of factual accuracy, evidence alignment, and risk minimisation in the evaluation rubric proved important in capturing these distinctions and avoiding the conflation of correctness with neutrality.

Performance differences across content categories further support this interpretation. Scientific queries, where evidentiary boundaries are clearer, yielded the highest agreement with experts, while marketing and regulatory queries posed greater challenges due to their reliance on framing and value-laden assumptions. Nevertheless, the framework remained effective across all categories, suggesting that the underlying evaluation logic generalises beyond narrowly defined factual questions.

**6.2 Contribution to Bias Evaluation in Public Health AI**

This work contributes to the growing literature on auditing LLM behaviour in high-stakes domains by demonstrating the feasibility of a structured, agent-based approach tailored to public health contexts. Unlike many existing bias evaluations that focus on demographic stereotypes or sentiment polarity, the present framework addresses domain-specific bias rooted in historical corporate influence and strategic communication.

A key contribution is the explicit separation between **evidence synthesis** and **bias evaluation**. By assigning these functions to distinct agents, the framework mirrors expert review practices and reduces the risk that evaluative judgements are driven by unconstrained, post-hoc reasoning. The use of an evidence-aligned baseline also provides a principled reference point for evaluation while remaining flexible to evolving scientific knowledge.

The structured output format further enhances transparency. Numeric scores allow comparison across models and queries, while accompanying qualitative annotations support interpretability and enable targeted error analysis. This dual-level output is particularly important in public health applications, where understanding *why* a response is problematic is often as important as identifying that it is.

## 6.3 Limitations

Several limitations should be acknowledged. First, although the evidence-aligned baseline is grounded in authoritative sources, it is generated through an LLM-guided synthesis process rather than derived from a fixed, human-curated gold standard. While expert validation suggests that this approach is generally reliable, it nonetheless introduces an element of epistemic dependence on LLM behaviour. The framework mitigates this risk through constrained retrieval, explicit documentation of uncertainty, and transparent reporting of baseline sources, but future work could incorporate hybrid human–AI baselines for particularly sensitive queries.

Second, bias assessment necessarily involves normative judgement, particularly when evaluating framing and omission. Although the rubric-based approach improves consistency, some degree of subjectivity remains unavoidable. Disagreements between the framework and experts highlight areas where evaluative criteria could be refined or made more explicit, rather than indicating clear errors.

Third, the evaluation focused on a curated set of tobacco-related queries. While these queries were designed to capture a range of known industry-aligned narratives, they do not exhaust the full space of possible prompts or contexts in which bias may arise. Generalisation to other domains should therefore be undertaken with care and domain-specific adaptation.

## 6.4 Implications and Future Work

The findings have practical implications for both researchers and practitioners deploying LLMs in public health contexts. The framework provides a scalable method for auditing model outputs prior to deployment, supporting risk assessment and model selection in settings where biased information could have real-world consequences. It also offers a tool for ongoing monitoring, enabling the detection of shifts in model behaviour as training data and deployment contexts evolve.

Future work could extend the framework in several directions. Incorporating multiple independent baseline generators could reduce dependence on a single synthesis pathway and enable inter-baseline agreement analysis. Expanding the rubric to capture additional dimensions, such as conflict-of-interest disclosure or regulatory framing, may further improve sensitivity. Finally, adapting the framework to other public health domains affected by corporate influence,

such as alcohol, ultra-processed foods, or pharmaceuticals, represents a promising avenue for broader application.

## 7. Conclusion

This study introduced and evaluated a multi-agent framework for detecting tobacco industry–favourable bias in large language model outputs. By separating evidence synthesis from evaluative judgement and implementing both steps within a unified CrewAI architecture, the framework provides a transparent and reproducible approach to auditing LLM behaviour in a high-risk public health domain.

The results demonstrate that a structured, rubric-based comparison between model responses and an evidence-aligned baseline can approximate expert judgement with a high degree of agreement. Importantly, the framework does not equate bias with factual error. Instead, it captures subtler forms of industry-aligned framing, such as risk minimisation, selective emphasis, and omission of key caveats, which are central to contemporary tobacco industry communication strategies.

Beyond tobacco control, the methodological contribution of this work lies in showing how agent-based systems can be used not to generate content, but to *evaluate* it in a principled and auditable manner. The framework's modular design, reliance on explicit scoring rubrics, and combination of quantitative and qualitative outputs make it adaptable to other public health and policy domains where corporate influence and strategic framing pose ongoing challenges.

While further work is needed to strengthen epistemic robustness and explore generalisation beyond the evaluated query set, the proposed framework represents a practical step toward systematic auditing of LLM outputs in contexts where biased information may have significant societal consequences. As LLMs continue to be integrated into research, policy, and public communication, such evaluation mechanisms will be essential for ensuring that technological advances do not inadvertently reproduce longstanding patterns of misinformation and influence.

The proposed framework provides a transparent, rubric-based approach for auditing industry-aligned bias in large language model outputs. By separating evidence grounding from bias evaluation and enabling systematic comparison across models, the approach supports scalable

auditing in high-risk public health domains where framing and omission can materially influence interpretation.

## Data and Code Availability

Source code, configuration files, and example datasets are publicly available at [repository URL].

## Competing Interests

The authors declare no competing financial interests. This research was conducted to advance public health and AI safety without industry funding or involvement.

## References

[1] The Lancet Digital Health. (2023). Generative AI: The future of digital health. *Lancet Digital Health*, 5(5), e262.

[2] Thirunavukarasu, A. J., et al. (2023). Large language models in medicine. *Nature Medicine*, 29(8), 1930-1940.

[3] Rani, P., et al. (2024). Bias and trustworthiness in artificial intelligence: A systematic review. *AI and Ethics*, 4(1), 123-145.

[4] Templin, M., et al. (2025). Algorithmic bias in healthcare applications of large language models. *Journal of Medical AI*, 12(2), 234-256.

[5] Ulucanlar, S., Fooks, G. J., & Gilmore, A. B. (2014). The policy dystopia model: An interpretive analysis of tobacco industry political activity. *PLoS Medicine*, 11(9), e1001738.

[6] Li, M., et al. (2025). Evaluating demographic bias in clinical large language models. *NPJ Digital Medicine*, 8(1), 45-58.

[7] Thapa, C., et al. (2024). Transformer-based language models for detecting health misinformation: A systematic review. *Journal of Biomedical Informatics*, 145, 104458.

[8] Wang, R., et al. (2025). Can large language models identify and explain health misinformation? *JMIR Medical Informatics*, 13(1), e45789.

[9] Yao, S., et al. (2023). ReAct: Synergizing reasoning and acting in language models. *International Conference on Learning Representations (ICLR)*.

[10] Yehudai, G., et al. (2025). Multi-agent collaboration patterns in complex reasoning tasks. *Artificial Intelligence*, 321, 103923.

[11] Rani, P., et al. (2024). Trustworthiness challenges in large language models: A comprehensive survey. *ACM Computing Surveys*, 56(4), 1-38.

[12] Templin, M., & Zhang, Y. (2025). Healthcare AI fairness: Addressing algorithmic bias in medical language models. *Journal of Medical Systems*, 49(2), 78-92.

[13] The Lancet Digital Health. (2023). The future of LLMs in clinical medicine. *Lancet Digital Health*, 5(5), e262-e263.

[14] Thirunavukarasu, A. J., et al. (2023). Systematic evaluation of large language models in clinical applications. *Nature Medicine*, 29(8), 1930-1940.

[15] Thapa, C., & Jang, S. (2024). Benchmarking LLMs for health misinformation detection. *IEEE Access*, 12, 34567-34582.

[16] Wang, R., et al. (2025). Evaluating LLM capabilities in health news quality assessment. *Digital Health*, 11, 20552076241234567.

[18] United States v. Philip Morris USA, Inc., 449 F. Supp. 2d 1 (D.D.C. 2006).

[19] Paranjape, B., et al. (2023). Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems (NeurIPS)*, 36.

[20] Hugging Face. (2024). Transformers Agents: Building advanced LLM applications. *Hugging Face Documentation*.

[21] CrewAI. (2024). CrewAI: Framework for orchestrating autonomous AI agents. Retrieved from https://crewai.com

[22] Significant Gravitas. (2023). AutoGPT: An experimental open-source attempt to make GPT-4 fully autonomous. *GitHub Repository*.

[24] Chen, L., et al. (2022). Machine learning classification of pro- and anti-vaping Twitter content. *Tobacco Control*, 31(4), 567-574.

[25] European Commission. (2024). The EU AI Act: Regulation on artificial intelligence. *Official Journal of the European Union*.