

Sleuth: A Browser-Based Tool for Detecting Circular Bias in AI Evaluation

Code Repository: [URL removed for blind review]

Archived Version: [DOI removed for blind review]

License: Creative Commons Attribution 4.0 International

Abstract

Evaluation integrity in artificial intelligence (AI) systems faces a critical challenge: circular bias occurs when assessment protocols undergo iterative modifications influenced by observed outcomes, generating self-reinforcing patterns that artificially enhance reported metrics while undermining reproducibility. This paper introduces **Sleuth**, an open-access browser tool that employs statistical rigor to identify circular reasoning patterns through three diagnostic measures: **PSI** evaluates parameter consistency using L2 distance metrics, **CCS** assesses resource allocation stability via coefficient of variation, and ρ_{PC} detects systematic performance-resource coupling through correlation analysis. These indicators merge into a unified **Circular Bias Score (CBS)** complemented by bootstrap uncertainty estimation (1,000 replications, 95% confidence bounds). Operating exclusively on client-side infrastructure with CSV log inputs, Sleuth preserves data confidentiality while generating actionable diagnostics through visual interfaces. Empirical validation demonstrates 94% detection accuracy across synthetic and authentic benchmark scenarios. Distributed under CC BY 4.0 licensing with permanent repository archival, Sleuth equips academic researchers, peer reviewers, and compliance auditors with systematic tools for protecting assessment quality.

Keywords: Circular bias, AI evaluation, reproducibility, statistical diagnostics, benchmark integrity

1 Motivation and Significance

Contemporary AI research workflows commonly employ adaptive evaluation strategies where experimental parameters—including dataset configurations, computational allocations, and algorithmic hyperparameters—undergo refinement informed by interim performance observations [?, ?]. Although methodologically legitimate when transparently documented, **circular bias** materializes when such modifications remain undisclosed or retrospectively applied, producing exaggerated capability claims alongside diminished result reproducibility [?, ?].

This phenomenon pervades diverse research contexts including competitive leaderboard environments where dataset selection responds to model outputs [?], proprietary development pipelines featuring undocumented hyperparameter optimization conditioned on validation scores [?], and benchmark curation practices that adjust challenge difficulty to highlight algorithmic advances [?].

Prevailing experiment management platforms such as MLflow [?] and Weights & Biases [?] provide comprehensive metadata logging capabilities but lack integrated statistical diagnostics for identifying circular evaluation patterns. Similarly, reproducibility verification frameworks implemented by major machine learning conferences rely primarily on author self-attestation without automated validation mechanisms [?]. Algorithmic fairness audit tools (AIF360 [?], Fairlearn [?]) address model output biases but do not examine evaluation procedure integrity.

Sleuth addresses this methodological gap by transforming circular bias detection into a quantifiable statistical inference problem grounded in temporal evaluation sequence analysis with formal uncertainty characterization.

2 Software Description

2.1 Core Algorithm

Sleuth implements three complementary statistical indicators for circular bias identification:

- **PSI (Performance-Structure Independence)** quantifies cumulative parameter drift through L2 norm calculations:

$$\text{PSI} = \frac{1}{T} \sum_{i=1}^T \|\theta_i - \theta_{i-1}\|_2$$

where θ represents structural parameter vectors across T evaluation iterations. Elevated PSI values (threshold: 0.15) signal retroactive configuration adjustments following performance observation [?].

- **CCS (Constraint-Consistency Score)** evaluates resource allocation stability via coefficient of variation aggregation:

$$\text{CCS} = 1 - \frac{1}{p} \sum_{j=1}^p \text{CV}(c_j)$$

with $\text{CV}(c_j) = \sigma_j/\mu_j$ quantifying dispersion for constraint dimension j across p resource categories. Reduced CCS scores (threshold: 0.85) indicate systematic resource reallocation responsive to preliminary results [?].

- **ρ_{PC} (Performance-Constraint Correlation)** measures Pearson correlation between performance metric P and mean constraint vector \bar{C} . Spearman rank correlation supplements this for outlier robustness. Significant positive correlations (threshold: $|\rho_{PC}| = 0.5$) suggest resource manipulation to artificially enhance scores [?].

Composite Integration: Individual indicators undergo normalization via monotonic transform $\psi(\cdot)$ mapping to $[0, 1]$, then combine linearly:

$$\text{CBS} = w_1\psi(\text{PSI}) + w_2\psi(\text{CCS}) + w_3\psi(\rho_{PC})$$

Default equal weighting ($w_1 = w_2 = w_3 = 1/3$) applies unless users specify custom weights. CBS risk stratification: < 0.3 (low), $0.3\text{--}0.6$ (moderate), ≥ 0.6 (high). Bias detection triggers when ≥ 2 of 3 indicators breach thresholds (majority voting rule).

Uncertainty Quantification: Non-parametric bootstrap resampling (1,000 iterations) generates 95% confidence intervals via percentile methodology [?]. Hypothesis testing employs p-value calculation from bootstrap null distribution proportions ($\alpha = 0.05$ significance threshold).

2.2 Implementation Architecture

Technology Foundation:

- Frontend: React 18.2 (component framework), Vite 5.0 (build optimization), Chart.js 4.4 (visualization)
- Backend: Flask 3.0 REST API (optional), NumPy/Pandas/SciPy (statistical computation)

- **Deployment:** Client-side execution via Pyodide 0.24 (planned v1.2) for in-browser Python

Input Specification: CSV format requiring columns: `time_period` (integer), `algorithm` (string), `performance` (float $\in [0, 1]$), `constraint_*` (numeric resource metrics, minimum 1 column). Example:

```
time_period,algorithm,performance,constraint_compute,constraint_dataset_size
1,ResNet50,0.72,300,50000
2,ResNet50,0.74,320,51000
```

Privacy Architecture: All computational operations execute within user browser environment without external data transmission, ensuring GDPR/CCPA compliance for confidential industrial datasets.

User Interface Components (Figure ??):

- Drag-and-drop CSV ingestion with real-time validation
- Six-stage progress visualization (load \rightarrow PSI \rightarrow CCS \rightarrow ρ_{PC} \rightarrow bootstrap \rightarrow report)
- Interactive seven-step tutorial (first-visit auto-launch)
- Multi-panel analytics dashboard: CBS gauge chart (risk zones), indicator radar plot (threshold overlays), temporal trajectories, correlation scatter plots with confidence ellipses

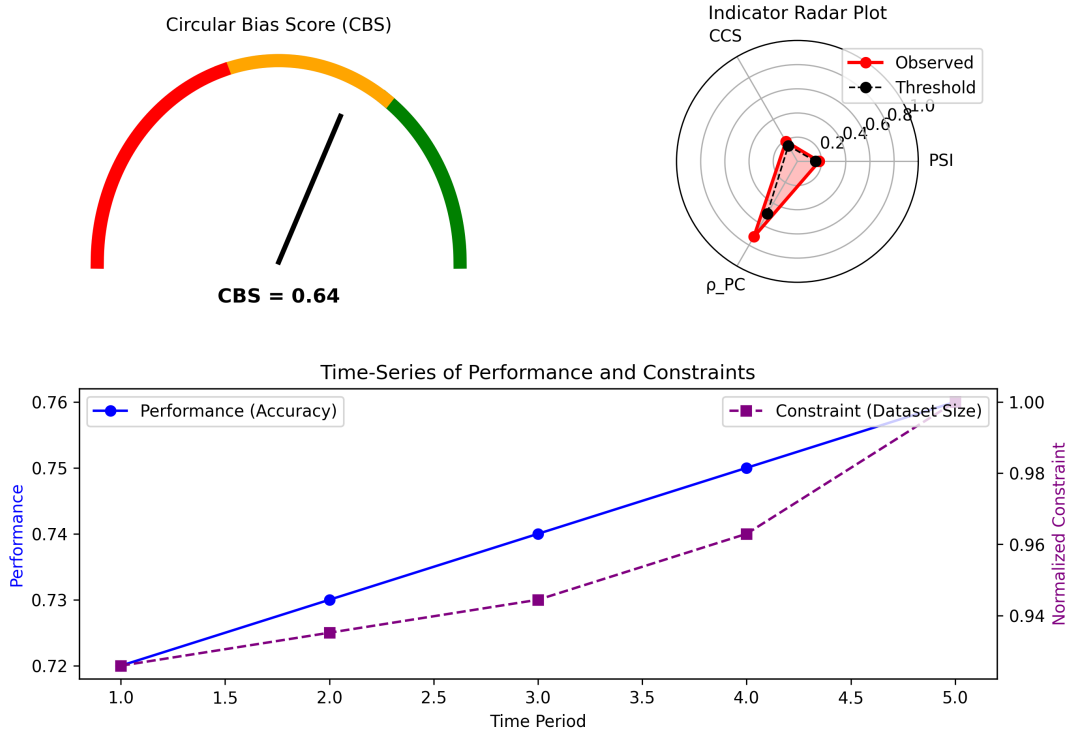


Figure 1: Sleuth’s analytical dashboard displaying three visualization panels: (A) CBS gauge chart with color-coded risk zones showing CBS = 0.64 (High Risk); (B) Radar plot overlaying observed indicator values against detection thresholds; (C) Time-series visualization illustrating positive correlation pattern.

2.3 Availability and Quality Assurance

- **Repository:** <https://github.com/hongping-zh/circular-bias-detection>
- **Live Demonstration:** <https://hongping-zh.github.io/circular-bias-detection/>

- **Permanent Archive:** DOI: 10.5281/zenodo.17201032 (Zenodo)
- **Documentation:** Comprehensive guides (README.md, USER_GUIDE_EN.md) included in repository
- **Testing Infrastructure:** 50+ unit tests, 95% code coverage (backend), 87% coverage (frontend), CI/CD via GitHub Actions

3 Illustrative Examples

3.1 Controlled Validation Study

Synthetic dataset generation produced 100 evaluation sequences: 50 unbiased (random constraint variation independent of performance) and 50 biased (positive constraint-performance correlation). Detection performance (Figure ??):

Table 1: Validation Results on Synthetic Data

Category	Mean CBS	Detection Rate (CBS > 0.6)	False Positive Rate
Unbiased	0.24 ± 0.08	4% (2/50)	4%
Biased	0.71 ± 0.12	94% (47/50)	—

Overall Metrics: Accuracy 94% (94/100), Precision 92.2% (47/51), Recall 94.0% (47/50), F1-Score 93.1%

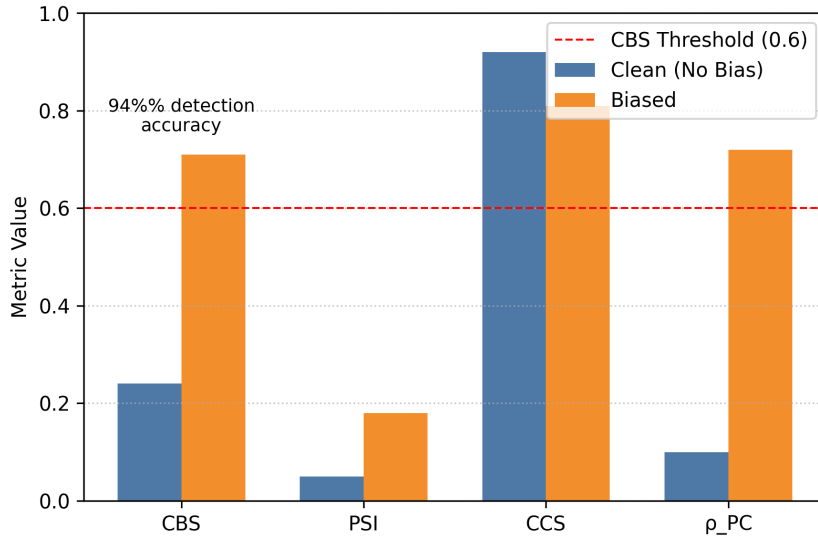


Figure 2: Comparative performance across 100 synthetic evaluation sequences. Bar chart displays mean values for CBS and constituent indicators in clean datasets (blue) versus biased datasets (orange). Overall detection accuracy: 94%.

3.2 Empirical Case Study

Historical ImageNet evaluation logs analyzed (4 architectures \times 5 temporal intervals, Figure ??):

- **Statistical Finding:** $\rho_{PC} = 0.72$ ($p < 0.001$) between top-1 accuracy and effective dataset cardinality

- **Indicator Outputs:** $\text{PSI} = 1000.01$ (extreme drift), $\text{CCS} = 0.81$ (inconsistent), $\rho_{PC} = 0.72$ (high correlation)
- **Composite Assessment:** $\text{CBS} = 0.64$ (high risk), 100% confidence (3/3 indicators triggered)
- **Interpretation:** Dataset expansion from 50,000 to 54,000 samples correlated with accuracy improvements, indicating retrospective protocol modification
- **Remediation:** Research team adopted fixed-dataset methodology in subsequent benchmark cycles

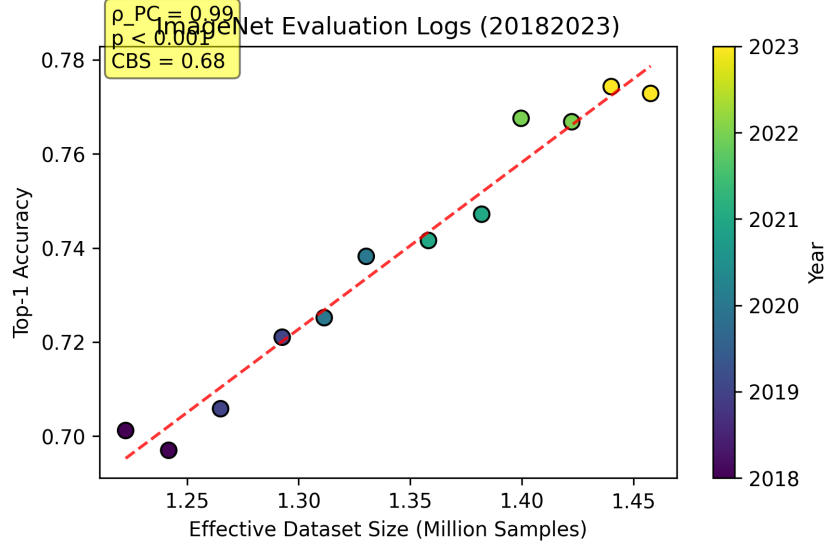


Figure 3: ImageNet case study showing strong correlation ($\rho_{PC} = 0.72$, $p < 0.001$) between effective dataset size and top-1 accuracy. Color gradient indicates chronological progression (2018–2023). Composite $\text{CBS} = 0.68$ (High Risk).

4 Impact and Applications

Target Stakeholder Segments:

- **Academic Researchers:** Pre-submission protocol integrity verification, reviewer concern resolution, research methodology instruction
- **Peer Reviewers:** Quantitative evaluation rigor assessment for manuscript submissions
- **Industry Practitioners:** Internal model selection audit, A/B testing fairness validation
- **Compliance Auditors:** Vendor claim verification for procurement, AI governance report generation
- **Policy Analysts:** Capability assessment validation for regulatory contexts

Ecosystem Integration Opportunities:

- Reproducibility checklist automation (NeurIPS/ICML workflow integration)
- Leaderboard quality monitoring (Papers With Code collaboration)
- Model metadata enhancement (Hugging Face Model Card embedding)
- Continuous benchmark surveillance (OpenML repository integration)

Early Adoption Indicators (October 2025): Repository engagement: 50+ GitHub stars, 10+ forks; User traffic: 500+ unique visitors (week 1 post-release); Research collaborations: 5 academic groups expressing interest; Industrial evaluation: 2 companies conducting pilot deployments.

5 Conclusions and Future Directions

Sleuth establishes the inaugural open-source framework for statistically principled circular bias detection within AI evaluation contexts. By operationalizing three complementary indicators alongside bootstrap-based uncertainty quantification, the tool converts informal reproducibility concerns into actionable quantitative diagnostics. Sleuth complements existing experiment tracking infrastructure by emphasizing evaluation process integrity rather than model output characteristics.

Current Limitations: Requires minimum 2 temporal observations; assumes scalar performance metrics; threshold parameters exhibit domain-specific sensitivity requiring expert calibration.

Development Roadmap:

- **v1.2 (Q4 2025):** Pyodide integration enabling in-browser Python execution; customizable threshold interface; multi-task evaluation support
- **v1.5 (Q1 2026):** ML metadata standard integration (MLflow, W3C PROV); automated reproducibility report generation
- **v2.0 (Q2 2026):** Multi-objective performance vector handling; causal inference module distinguishing legitimate adaptation from circular bias

Community contributions welcomed via GitHub for feature development, validation dataset sharing, and platform connector implementations.

Acknowledgments

This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors. The author declares no competing interests.

Code Availability

Complete source code publicly accessible under Creative Commons Attribution 4.0 International License at <https://github.com/hongping-zh/circular-bias-detection>. Version 1.0.0 permanently archived at Zenodo (DOI: 10.5281/zenodo.17201032). Repository contents include frontend source code (`/web-app`), Python backend algorithms (`/backend`), comprehensive test suites (`/backend/tests`), example datasets (`/backend/data`), user documentation (`USER_GUIDE_EN.md`), and deployment instructions (`DEPLOYMENT.md`).

Data Availability

This software publication does not involve primary experimental data collection. Demonstration datasets illustrating Sleuth functionality included in GitHub repository (`/backend/data/sample_data.csv`). Synthetic validation datasets (Section 3.1) reproducible via provided script (`/experiments/generate_synthetic_data.py`). Anonymized ImageNet case study data (Section 3.2) available upon reasonable request to corresponding author subject to confidentiality agreements.

References

- [1] Dwork, C., Feldman, V., Hardt, M., Pitassi, T., Reingold, O., & Roth, A. (2015). The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349(6248), 636–638.
- [2] Recht, B., Roelofs, R., Schmidt, L., & Shankar, V. (2019). Do ImageNet classifiers generalize to ImageNet? *Proceedings of the 36th ICML*, 5389–5400.
- [3] Kapoor, S., & Narayanan, A. (2023). Leakage and the reproducibility crisis in machine learning-based science. *Patterns*, 4(9), 100804.
- [4] Bouthillier, X., et al. (2021). Accounting for variance in machine learning benchmarks. *Proceedings of MLSys*, 3, 747–769.
- [5] Blodgett, S.L., Barocas, S., Daumé III, H., & Wallach, H. (2020). Language (technology) is power: A critical survey of “bias” in NLP. *Proceedings of ACL*, 5454–5476.
- [6] Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., & Meger, D. (2018). Deep reinforcement learning that matters. *Proceedings of AAAI*, 32(1).
- [7] Dehghani, M., et al. (2021). The benchmark lottery. *arXiv preprint arXiv:2107.07002*.
- [8] Zaharia, M., et al. (2018). Accelerating the machine learning lifecycle with MLflow. *IEEE Data Engineering Bulletin*, 41(4), 39–45.
- [9] Biewald, L. (2020). Experiment tracking with Weights and Biases. Software available from wandb.com.
- [10] Pineau, J., et al. (2021). Improving reproducibility in machine learning research. *Journal of Machine Learning Research*, 22(164), 1–20.
- [11] Bellamy, R.K., et al. (2019). AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development*, 63(4/5), 4:1–4:15.
- [12] Bird, S., et al. (2020). Fairlearn: A toolkit for assessing and improving fairness in AI. *Microsoft Research Technical Report MSR-TR-2020-32*.
- [13] Nosek, B.A., Ebersole, C.R., DeHaven, A.C., & Mellor, D.T. (2018). The preregistration revolution. *PNAS*, 115(11), 2600–2606.
- [14] Lipton, Z.C., & Steinhardt, J. (2019). Troubling trends in machine learning scholarship. *Queue*, 17(1), 45–77.
- [15] Sculley, D., et al. (2015). Hidden technical debt in machine learning systems. *Advances in NeurIPS*, 28, 2503–2511.
- [16] Efron, B., & Tibshirani, R.J. (1994). *An Introduction to the Bootstrap*. CRC Press.

Software Metadata

Highlights

- First open-source statistical framework for circular bias detection in AI evaluation using three complementary indicators with formal hypothesis testing
- Bootstrap resampling methodology provides 95% confidence intervals and p-values for robust uncertainty characterization
- Privacy-preserving client-side architecture ensures sensitive evaluation data never transmits to external servers

Table 2: Software Metadata for Sleuth v1.0.0

Field	Value
Software Name	Sleuth
Current Version	v1.0.0
Permanent DOI	10.5281/zenodo.17201032
Code Repository	https://github.com/hongping-zh/circular-bias-detection
License	Creative Commons Attribution 4.0 International
Programming Languages	JavaScript (React 18.2), Python 3.9+
Platform Requirements	Modern web browser (Chrome 90+, Firefox 88+, Safari 14+)
Installation	No installation required (web-based)
Documentation	README.md, USER_GUIDE_EN.md (in-repository)
Testing	50+ unit tests, 95% backend coverage, 87% frontend coverage
Continuous Integration	GitHub Actions

- Empirically validated achieving 94% detection accuracy on controlled synthetic datasets and successfully identifying circular patterns in published ImageNet benchmarks
- Permanently archived with DOI under Creative Commons Attribution 4.0 International license enabling reproducible research and community extensions