# Machine Learning Techniques for Detecting Greenwashing in Corporate Financial Reports

*Gianna Foster, Grace Brooks, Hailey Butler*

**Abstract**

This research introduces a novel, hybrid machine learning framework specifically designed to detect and quantify greenwashing—the practice of making misleading environmental claims—within the narrative sections of corporate financial reports. While existing literature primarily focuses on sentiment analysis or keyword spotting for sustainability reporting, our approach uniquely integrates three distinct methodologies to capture the nuanced, often obfuscated nature of greenwashing. First, we employ a transformer-based language model fine-tuned on a purpose-built corpus of verified greenwashing and legitimate sustainability disclosures to perform deep semantic analysis, moving beyond surface-level features. Second, we implement a novel coherence scoring mechanism that measures the alignment between environmental claims made in the front-of-report narratives and the quantitative environmental performance data presented in appendices or supplementary reports, identifying strategic decoupling. Third, we develop a temporal inconsistency detector using recurrent neural networks to flag claims that contradict a company's own historical environmental disclosures. We validate our framework on a manually annotated dataset of 500 annual reports from the SP 500 between 1995 and 2004, achieving a detection accuracy of 91.7% and a precision of 88.3% in identifying materially misleading statements, significantly outperforming baseline keyword-matching and sentiment analysis models. Our findings reveal that greenwashing is not merely a function of exaggerated positive sentiment but is characterized by specific rhetorical patterns, strategic vagueness, and measurable disconnects between narrative and data. This work provides auditors, regulators, and investors with a powerful, automated tool for enhanced scrutiny of corporate environmental communications and establishes a new methodological paradigm for computational analysis of corporate discourse.

**Keywords:** Greenwashing Detection, Natural Language Processing, Corporate Reporting, Financial Disclosures, Machine Learning, Sustainability

# 1 Introduction

The proliferation of corporate sustainability reporting over the past two decades has been accompanied by a growing concern over greenwashing, wherein organizations exaggerate or misrepresent their environmental performance to cultivate a positive public image. This practice undermines the integrity of environmental, social, and governance (ESG) metrics, misleads stakeholders, and poses significant risks to sustainable investment. Traditionally, the detection of greenwashing has relied on manual analysis by domain experts, a process that is time-consuming, subjective, and difficult to scale across the vast corpus of corporate publications. While computational linguistics and machine learning have been applied to financial sentiment analysis and fraud detection, their application to identifying deceptive environmental rhetoric remains nascent and methodologically limited. Existing approaches often rely on simplistic keyword dictionaries or sentiment polarity scores, which fail to capture the sophisticated linguistic strategies employed in corporate greenwashing, such as the use of vague aspirational language, selective disclosure of positive information, or narrative-data decoupling.

This paper addresses this critical gap by proposing and validating a novel, multi-faceted machine learning framework for the automated detection of greenwashing in the narrative sections of annual financial reports. Our research is guided by two primary questions: First, can a hybrid machine learning model, integrating semantic, coherence, and temporal analysis, reliably identify instances of greenwashing with higher accuracy than existing keyword or sentiment-based methods? Second, what are the distinctive linguistic and rhetorical features that characterize greenwashing in formal corporate disclosures, as opposed to merely positive environmental messaging? Our contribution is threefold. Methodologically, we introduce a new paradigm that moves beyond bag-of-words models to analyze the deeper semantic structure, internal consistency, and historical fidelity of environmental claims. Empirically, we present a new, manually annotated dataset of corporate reports labeled for greenwashing, serving as a benchmark for future research. Practically, we deliver a tool that can enhance the monitoring and enforcement capabilities of regulators, the due diligence processes of investors, and the assurance practices of

auditors.

# 2    Methodology

Our proposed framework, the Greenwashing Detection Integrated Model (GDIM), consists of three interconnected analytical modules designed to operate on the full text of a corporate annual report. The input is segmented into narrative sections (e.g., CEO letter, management discussion) and quantitative appendices containing environmental metrics.

The first module is the Semantic Deception Classifier (SDC). Instead of using pretrained general-purpose language models, we constructed a specialized training corpus. This corpus comprises text snippets from two sources: confirmed cases of greenwashing as adjudicated by regulatory bodies like the U.S. Federal Trade Commission and the U.K. Advertising Standards Authority between 1998 and 2004, and verified, substantive environmental disclosures from sustainability reports certified by third-party auditors. A fine-tuned transformer model, building on the architectural principles of attention mechanisms, learns to distinguish the nuanced language of deception from that of legitimate reporting. It identifies patterns such as excessive use of vague, non-actionable terms (e.g., "committed to," "aiming for"), disproportionate focus on minor environmental initiatives while omitting major impacts, and the use of emotional appeals divorced from concrete plans.

The second module is the Narrative-Data Coherence Scorer (NDCS). This novel component addresses a core greenwashing tactic: making bold claims in the narrative that are unsupported or contradicted by the hard data presented elsewhere in the report. The module first uses a rule-based information extraction system to identify quantifiable environmental claims in the narrative (e.g., "reduced emissions by 20%"). It then locates corresponding data points in tables, charts, or footnotes. A coherence score is calculated based on the logical alignment between the claim and the data. A significant negative score triggers a greenwashing flag. For instance, a narrative claiming "significant investment in renewable energy" paired with data showing a decrease in renewable energy

3

spending would yield a low coherence score.

The third module is the Temporal Inconsistency Detector (TID). Greenwashing often involves claims that represent a break from a company's historical trajectory or previous commitments without adequate explanation. This module employs a recurrent neural network (RNN) architecture to model a company's sequence of annual reports. It learns the expected progression of environmental discourse and performance for a given sector. When a new report contains claims that represent a statistically significant positive deviation from the learned historical pattern without a corresponding deviation in underlying performance data, it is flagged for potential greenwashing. This captures instances of sudden, unsubstantiated "green rebranding."

The outputs of these three modules—a deception probability from the SDC, a coherence score from the NDCS, and an inconsistency flag from the TID—are fused using a weighted meta-classifier (a support vector machine) to produce a final, holistic greenwashing likelihood score for the document. The weights for the meta-classifier were optimized on a validation set.

# 3  Results

We evaluated the GDIM framework on a curated dataset of 500 annual reports from companies listed on the SP 500 index, spanning the decade from 1995 to 2004. This period saw a rapid increase in voluntary environmental reporting but preceded modern ESG standardization, making it a rich environment for studying greenwashing. Each report was manually annotated by a panel of three experts in environmental accounting and corporate communication. Annotations identified specific sentences or paragraphs as instances of clear greenwashing, legitimate positive disclosure, or neutral information. Inter-annotator agreement, measured by Fleiss' kappa, was 0.81, indicating substantial reliability.

The performance of the complete GDIM was compared against two strong baseline models: a keyword-matching model using an extensive dictionary of "green" terms cou-

pled with sentiment analysis, and a logistic regression model using standard TF-IDF (Term Frequency-Inverse Document Frequency) features. The results demonstrate the superiority of our integrated approach. The GDIM achieved an overall accuracy of 91.7% and a precision of 88.3% for the greenwashing class, meaning that when it flagged a statement as greenwashing, it was correct 88.3% of the time. Its recall was 85.9%. In contrast, the keyword-sentiment baseline achieved an accuracy of 72.1% and a precision of only 61.5%, frequently misclassifying genuinely positive but enthusiastically worded disclosures as greenwashing. The TF-IDF model performed slightly better with 78.3% accuracy but lacked interpretability.

Ablation studies, where individual modules of the GDIM were disabled, revealed the contribution of each component. Using only the Semantic Deception Classifier (SDC) yielded an accuracy of 84.2%. Adding the Coherence Scorer (NDCS) boosted accuracy to 89.1%, and the full model with the Temporal Detector (TID) reached the peak performance of 91.7%. This confirms that each module captures a distinct and complementary aspect of greenwashing rhetoric.

Qualitative analysis of the model's outputs provided novel insights into the nature of greenwashing. We found that the most reliable indicators were not simply positive sentiment, but specific linguistic constructs: the use of the future tense to defer accountability ("we will become carbon neutral"), nominalizations that obscure agency ("a reduction was achieved"), and what we term "contextual isolation"—highlighting a single green product line while the company's core business remains environmentally intensive. The NDCS module revealed that over 40% of reports with high greenwashing scores contained at least one major narrative claim that was directly contradicted by data in the same report, often buried in technical footnotes.

# 4    Conclusion

This research has established that machine learning techniques, when designed to address the specific rhetorical and structural complexities of corporate environmental disclosures,

can effectively automate the detection of greenwashing. Our proposed Greenwashing Detection Integrated Model (GDIM) represents a significant departure from and improvement over previous computational methods, which were largely inadequate for this subtle task. By integrating deep semantic understanding, cross-document coherence analysis, and temporal pattern recognition, the GDIM achieves high accuracy in distinguishing between legitimate sustainability communication and deceptive greenwashing.

The original contributions of this work are manifold. Methodologically, we have introduced a new framework that combines multiple NLP and machine learning paradigms in a novel way for a novel problem. We have demonstrated the critical importance of moving beyond isolated text analysis to consider the relationship between narrative and data, and between present and past disclosures. Empirically, we have generated new knowledge about the linguistic signatures of greenwashing, identifying patterns of vagueness, temporal displacement, and contextual isolation as key markers. The creation of our annotated dataset provides a valuable resource for the research community.

The practical implications are substantial. For regulators such as the Securities and Exchange Commission, this tool could enable large-scale, continuous monitoring of corporate reports for misleading claims. For asset managers and ESG rating agencies, it offers an objective, scalable layer of analysis to complement human judgment, potentially reducing "greenwashing risk" in investment portfolios. For companies themselves, it could serve as a self-assessment tool to improve the integrity of their communications.

Future work will focus on expanding the model to analyze multimedia corporate communications (e.g., sustainability videos, website content) and to incorporate sector-specific linguistic models, as greenwashing tactics may vary between industries such as energy, manufacturing, and finance. Furthermore, the framework could be adapted to detect analogous phenomena like "social washing" or "governance washing." In conclusion, this research bridges a critical gap between computational linguistics and sustainable finance, providing a sophisticated, evidence-based tool to promote transparency and accountability in corporate environmental reporting.

# References

Bowen, F. E. (2002). Organizational slack and corporate greening: Broadening the debate. *British Journal of Management, 13*(4), 305-316.

Cho, C. H., Patten, D. M. (2004). Green accounting: Reflections from a CSR and environmental disclosure perspective. *Critical Perspectives on Accounting, 15*(6-7), 731-741.

Deegan, C. (2002). The legitimising effect of social and environmental disclosures–a theoretical foundation. *Accounting, Auditing & Accountability Journal, 15*(3), 282-311.

Feldman, R., Sanger, J. (2002). *The text mining handbook: Advanced approaches in analyzing unstructured data.* Cambridge University Press.

Joachims, T. (1998). Text categorization with support vector machines: Learning with many relevant features. In *European conference on machine learning* (pp. 137-142). Springer.

Laufer, W. S. (2003). Social accountability and corporate greenwashing. *Journal of Business Ethics, 43*(3), 253-261.

Loughran, T., McDonald, B. (2004). What's in a word? A study of the tone of corporate financial disclosures. *Journal of Finance, 59*(3), 1457-1489.

Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., Dean, J. (2003). Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems, 26.*

Ramus, C. A., Montiel, I. (2005). When are corporate environmental policies a form of greenwashing? *Business & Society, 44*(4), 377-414.

Waddock, S. A., Graves, S. B. (1997). The corporate social performance-financial performance link. *Strategic management journal, 18*(4), 303-319.