

**ENHANCING THE CATEGORIZATION OF SUSTAINABILITY REPORTS  
USING A MULTIMODAL DEEP LEARNING APPROACH**

**SARAH CHEPKOGEI SAWE**

**A PROJECT SUBMITTED TO THE DEPARTMENT OF COMPUTER SCIENCE  
& INFORMATION TECHNOLOGY IN THE SCHOOL OF COMPUTING AND  
MATHEMATICS FOR PARTIAL FULFILLMENT OF THE REQUIREMENT OF  
THE AWARD OF MASTER OF SCIENCE IN DATA SCIENCE OF THE  
COOPERATIVE UNIVERSITY OF KENYA**

**2025**

## DECLARATION

### Declaration by the Candidate

This project is my original work and has not been presented for award of a degree in any other University or for any other award.

.....

Signature

.....

Date

Sarah Chepkogei Sawe

MDATC01/6062/2022

### Declaration by the supervisors

We confirm that the work reported in this project was carried out by the candidate under our supervision and has been submitted with our approval as university supervisors.

.....

Signature

.....

Date

Dr. Anthony Wanjoya

Department of Computer Science & Information Technology

The Cooperative University of Kenya

.....

Signature

.....

Date

Dr. Andrew Kipkebut

Department of Computer Science & Information Technology

The Kabarak University

## **DEDICATION**

To all who promote sustainable growth and new solutions for the future. Their drive, enthusiasm, and resilience motivate them to make lasting change. May this work help them

## **ACKNOWLEDGEMENT**

I would like to begin by Thanking God Almighty for the grace, strength and opportunity to complete my research.

My sincere gratitude to the Cooperative University of Kenya for giving me opportunity to pursue my master's degree and for providing a supportive environment to complete this research.

This research is the result of a joint effort, and I'd like to thank my supervisors, Dr. Anthony Wanjoya, Dr. Ronald Ojino, and Dr. Andrew Kipkebut, for their invaluable time, experience, direction, and support throughout the journey.

Special thanks to my family and classmates for their encouragement, insights, and continuous support in making this work possible.

## TABLE OF CONTENTS

DECLARATION .....	i
DEDICATION .....	ii
ACKNOWLEDGEMENT .....	iii
TABLE OF CONTENTS.....	iv
LIST OF TABLES .....	vii
LIST OF FIGURES .....	viii
ABBREVIATIONS .....	ix
ABSTRACT.....	xii
CHAPTER ONE .....	1
1. INTRODUCTION .....	1
1.1. Background of study .....	1
1.2. Statement of the problem .....	8
1.3. Objectives of the study.....	9
1.3.1. General Objective .....	9
1.3.2. Specific objectives .....	9
1.4. Research Questions .....	9
1.5. Significance of the study.....	10
1.6. Scope of the study .....	11
1.7. Limitations of the study .....	11
CHAPTER TWO .....	12
2. LITERATURE REVIEW .....	12
2.1. Introduction.....	12
2.1.1. Theoretical Framework.....	12
2.1.2. General Literature Review .....	15
2.2. Conceptual Framework.....	21
2.3. Empirical Literature Review.....	23
2.4. Summary of Literature .....	25
2.5. Research gap .....	28
CHAPTER THREE .....	30

3. METHODOLOGY .....	30
3.1. Introduction.....	30
3.2. Research Philosophy.....	30
3.3. Research Design.....	31
3.3.1. Problem Identification .....	31
3.3.2. Define Objectives of a Solution.....	32
3.3.3. Design and Development of the Multimodal Deep Learning Model.....	32
3.3.4. Data Extraction and Model Training .....	35
3.4. Study Area .....	37
3.5. Target Area .....	38
3.6. Sampling Method.....	38
3.7. Data collection .....	40
3.8. Data Collection Procedures.....	42
3.9. Data analysis and presentation.....	43
3.10. Empirical Model and Hypothesis Testing.....	44
3.11. Ethical Considerations .....	46
CHAPTER FOUR.....	47
4. RESULTS AND DISCUSSION.....	47
4.1. Introduction.....	47
4.2. Addressing Research Objectives.....	47
4.2.1. Analysis of Multimodal Integration Impact.....	47
4.2.2. Design and Implementation of the Multimodal Framework.....	48
4.2.3. Development and Optimization of the Ensemble Model.....	49
4.3. Model Performance Evaluation. ....	49
4.3.1. Overall Classification Accuracy .....	49
4.3.2. Performance Across ESG Dimensions .....	51
4.3.3. Ablation Study Findings .....	52
4.3.4. Efficiency and Scalability Analysis .....	52
4.3.5. Addressing Research Questions.....	52
4.4. Model Performance Evaluation .....	53
CHAPTER FIVE .....	57

5. CONCLUSION AND RECOMMENDATIONS .....	57
5.1. Discussion and Finding.....	57
5.2. Conclusions.....	59
5.3. Recommendations.....	59
5.3.1. For Future Research.....	59
5.3.2. For Industry and Policy Makers.....	60
REFERENCES .....	61
Appendix A. University Approval Letter .....	66
Appendix B: Research Permit.....	67
Appendix C. Journal Publishing Certificate .....	68
Appendix D. Similarity Index Report.....	69
Appendix D. AI Content Percentage .....	71
Appendix E: Publication .....	73

## LIST OF TABLES

Table 2.1 Summary and Gaps in the Literature review ..... 25

Table 3.1 1: Components of the model under development ..... 33

## LIST OF FIGURES

Figure 2. 1 Conceptual Framework Source:(Author 2025) .....	22
---	----

## ABBREVIATIONS

AI:	Artificial Intelligence
BERT:	Bidirectional Encoder Representations from Transformers
CNNs:	Convolutional Neural Networks
CRR:	Corporate Responsibility Reporting
CSV:	Comma-Separated Values
CSR:	Corporate Social Responsibility
ESG:	Environmental, Social, and Governance
GRI:	Global Reporting Initiative
HTML:	HyperText Markup Language
KPIs:	Key Performance Indicators
KNN:	K-Nearest Neighbors
LSTMs:	Long Short-Term Memory Networks
MySQL:	My Structured Query Language
NLP:	Natural Language Processing
OCR:	Optical Character Recognition
PDF:	Portable Document Format
RNNs:	Recurrent Neural Networks
RoBERTa:	Robustly Optimized BERT Approach
SASB:	Sustainability Accounting Standards Board
SDGs:	Sustainable Development Goals
SMEs:	Small and Medium Enterprises
SVMs:	Support Vector Machines

UN: United Nations  
XML: Extensible Markup Language  
FCNs: Fully Connected Neural Networks

## **OPERATIONAL DEFINITIONS**

- Sustainability:** The ability to satisfy current needs, ranging from environmental, social, and economic factors without endangering the capacity of future generations to satisfy their own
- TensorFlow:** A popular tool for large-scale machine learning applications and numerical computing developed as an open-source machine learning framework by Google to create and train deep learning models.
- Deep Learning:** A branch of machine learning that makes use of multi-layered artificial neural networks to analyse and comprehend intricate data patterns; frequently used for image recognition, natural language processing, and other sophisticated AI applications.
- Multimodal:** A system or model capable of processing and integrating input from multiple modalities or data kinds, including text, images, audio, and video, to generate more comprehensive insights

## ABSTRACT

There is a growing need for organisations to align their operations and services with the sustainable development goals. In Kenya and most countries across the globe, it is mandatory requirement, for organisations especially those trading in the stock exchange, to provide their sustainability reports annually. The categorisation and classification process, however, for the reports has become very complex and demanding especially given the growing number of reports and the diverse nature of the data and information contained in them. Traditional categorisation process has proven inadequate, inefficient and not very accurate as compared to the modern process which integrates artificial intelligence and multimodal deep learning. Due to the aforementioned, this study aimed at developing an artificial intelligent model that can read various data formats, including textual, graphical and numeric and ensure capture the varied and intricate facts with accuracy and precision. The study further sort to conduct an analysis on the integration of the multimodal data from sustainability reports and understand its impacts on categorization accuracy and provides deeper ESG insights. The findings of the study informed the development of an AI model that automated the categorisation of the sustainability reports through the employing the Transformer Models, Recurrent Neural Networks (RNN), Convolutional Neural Networks (CNN), and Fully Connected Neural Networks (FCNN), in conjunction with high-performance computing resources. Approximately, 50 sustainability reports drawn from the submissions made by the organisations that are trading in Nairobi stock exchange and the Global Reporting Initiative database, were used to test the model. The results confirmed the model's robust performance, achieving an overall accuracy of 79.0%, an F1-score of 0.76, precision of 0.82, and a recall of 0.77. This performance demonstrated a significant enhancement in efficacy and accuracy compared to traditional unimodal baselines, thereby addressing the issue of inefficiency and cost of manual analysis. The model was also enhanced to extract key metadata, check publication dates, title company, and country, and provide percentages of alignment of the organizations' processes to the Sustainable Development Goals (SDGs). The model was designed to accommodate different types of deep learning models, thus can analyse images, texts and numerical. The findings indicate the significance of multimodal learning and advanced deep learning models in enhancing ESG reporting quality and competence to make evidence-based judgments in the sustainability sector.

## **CHAPTER ONE**

### **1. INTRODUCTION**

#### **1.1. Background of study**

Sustainability is at the center of every development agenda and a very important subject. Every aspect of humanity and its existence is affected by issues pertaining to sustainable development. Sustainable development involves achieving growth and human development by exploitation of the existing resources to meet current needs without disadvantaging the sustenance of future generations (Ozili, 2022).

Organizations across all fields are under immense pressure to ensure transparency, especially on their environmental, social and governance structures, and processes. This pressure, especially from the environmental sustainability perspective, comes from the UN member countries to uphold sustainable development in a way that meets the current societal human needs while making sure to preserve planetary integrity.

The increasing demand and emphasis on sustainability across the diverse industries has prompted a huge rise in the number of sustainability reports that are published. The importance of these documents cannot be understated as they are critical in offering a glimpse on organizations' sustainability practices which are essential for decision making, regulatory requirements amongst meeting other reporting mandates (Smith & Green, 2021). The sheer volume and diversity of these reports, which often contain a blend of textual, visual, and numerical data, severely hamper their effective categorization and analysis (Johnson et al., 2020). Traditionally, analysis and categorization of these documents over-relied on manual approaches such as rule-based algorithms which have

proven quite redundant, time consuming and inefficient, especially due to their limited ability to capture the full scope of multimodal data contained in a report.

Artificial Intelligence plays a critical role in mitigation of the challenges and gaps that lie within the conventional categorization models. Deep learning technology is at the center of innovations that is creating more suitable tools for managing big multidata (Zhang & Wang, 2018). Neural network- based deep learning models, continue to deliver key benefits in their areas of operations especially through classification of data, computer visions and natural language processing (NLP) (Brown & Harris, 2020). Classifications and categorisation that are based on multimodal approaches has proven to be more efficient, accurate and precise. This is due to the thorough integration of the various technology in the analysis of ESG related data within the multimodal work frames (Kim et al., 2022). A critical exploration of multimodal deep learning approaches and their great analytical potential is quite necessary and timely.

Bal et al., (2013) affirms that customers are no longer ignorant and are more attracted to organisations that possess and embrace sustainable behaviours and practices. Through the ESG reports, organisations can communicate to a wide range of stakeholders, existing and potential customers, and reaffirm their commitment to the sustainable environmental, social and governance structures. But while this stands, the ESG reports are increasingly becoming more detailed, complex and the reporting mandates are always increasing. This is coupled by the fact that reporting is a process and probably not the core organisational mandate, thus needn't be an expensive and labour-intensive process. The detailed nature of sustainability reports makes them complex and lengthy, as companies detail Environmental, Social, and Governance (ESG) metrics. In most cases, the process of

analyzing sustainability reports is done manually, which creates multiple issues, considering manual analysis is time-consuming and costly. Manual analysis harbors scalability, accuracy, and efficiency issues. Additionally, the manual process also comes with the risk of potential errors and missed points. The ever-evolving digital workplace presents an opportunity to solve the reports' analysis issue.

Rule based algorithm, which predominantly incorporates the 'if then' approach in its execution of problems, is more predictable and deterministic in character and will only produce results within the predefined environment. Durkin (1994) noted that rule-based algorithm usefulness is limited to only situations where transparency is necessary and has had immense impacts in running diagnostics in medicine and in the finance sector. Durkin (1994) further poses that, it is in the static nature of the rule-based algorithm that make them less adaptable over time. Crucially in the ESG domain, rule-based systems struggle with the inherent inconsistencies and non-standardization across different reporting frameworks (GRI, SASB, etc.) and fail to adequately process the complex, contextual language and diverse data formats (tables, charts, images) found in modern sustainability reports. Rule-based systems are also limited for lack of scalability, high maintenance costs and inefficiency with large data and major performance bottlenecks. Russel & Norvig, (2021), however, observed that machine learning models stand to be more beneficial over the rule-based algorithms due to their ability to refine process based on incorporating patterns and new data. Witten et al., (2016) also delved into the subject and highlighted the ability of email filtering systems using rule-based algorithms to classify emails and filter spam. He, however, noted that phishing scams have also evolved with time to become

more elusive to catch as opposed to the traditional phishing scam where emails contained phrases like ‘click to win money’ and were easy to mark.

Armed with speed and accuracy in analyzing large amounts of data, AI tools can help alleviate challenges of manual analysis of sustainability reports. This is because AI tools can analyze huge datasets and produce meaningful and timely insights, trends, and recommendations. According to Najafabadi et al., (2015), AI, particularly deep learning, improves data analysis in complex projects, which makes it a viable tool to eliminate issues with manual analysis of ESG reports, including time consumption and being prone to errors. This is because AI tools automate feature learning, recognize intricate patterns, and support transfer learning, which shows its enriched analysis capabilities (Najafabadi et al., 2015).

Data in ESG reports is presented in text, images, tables, and videos. It is necessary to have tools that have the capability to analyse such data i.e. in its different types and forms, concurrently and produce insights and calls for action. This necessitates a multimodal approach, which combines text, audio, and image as individual models. By combining data from multiple modalities, multimodal models can make more accurate and comprehensive predictions (Denissen & Chen , 2021).

Countries with access to more resources and research are better equipped to have such multimodal AI tools to aid in analyzing ESG reports than less developed nations. For example, in the USA where there is increased pressure for companies to align with sustainable development initiatives, firms have adopted AI tools, specifically deep learning models as major facilitators of reports’ analysis. For instance, Truvalue Labs uses AI to assess ESG data from multiple sources and can analyze information from over 100,000

sources (Tang, 2024). Truvalue uses AI to conduct analysis on big data in real-time. It is very important to take note of Truvalue's use of Natural Language Processing (NLP) to comprehend the relationship, frequency, and significance of ESG issues from multiple sources, including news, articles, blogs, as well as corporate reports (Tang, 2024). As a result, the AI algorithm helps stakeholders have updated insights on organizations' sustainability efforts, which drive informed decision-making.

RepRisk which can arguably be classified as a global leader in ESG technology, identifies and assesses ESG risks and is a leader in the provision of advanced machine learning analysis and research. The platform is capable of running among the world's largest, updated database in ESG risks daily. It uses multiple data sources, such as media and NGOs to produce risk factors, such as reputational risks. ESG risk is analyzed by combining AI and advanced machine learning with human intelligence (Lahey, 2024). The Truvalue Labs and RepRisk platforms prove that AI can rapidly and accurately aid in the process of analyzing ESG reports even when faced with unstructured data or data in multiple languages. While both platforms are powerful tools for ESG risk and sustainability analysis, their reliance on public data, potential biases, and cost considerations may limit their applicability in certain contexts, necessitating for a reliable open-source application that is multimodal with a broad scope of coverage and one that allows customization.

The implementation of AI milestones in global sustainability reporting encounters numerous challenges, notably the absence of high-quality datasets for training AI models. The capacity of AI to learn and reproduce biases included in the training data results in inequitable or discriminatory consequences. (Emilio, 2024). The potential for bias and data quality issues in AI cannot be underestimated, mainly since AI results depend on the data

fed into them, which means inaccurate data leads to skewed analysis and recommendations. Similarly, algorithm design issues and errors in human interpretation can also flaw results. A major noticeable issue at the global level is the integration of multimodal data. The challenge results from AI's ability to summarize data from many modalities effectively so that only crucial data is acquired, and the redundant part of the data is filtered out (Chen et al., 2024). At the same time, progress on globally accepted standards for AI-driven sustainability is still slow (Bashir, et al., 2024).

These global issues are amplified when it comes to AI's application in Africa. This is evidenced by African countries reduced efforts on sustainability matters, which can be explained by multiple barriers, including financial limitations, lack of knowledge and technology, and current policies and regulations (Setyaningsih et al., 2024). As a result, Africa trails other continents when it comes to using AI tools to analyze sustainability reports. However, progress is picking up speed, especially with the pressure on multinational companies operating in the region. Attention to sustainability reporting in Africa is also driven by rapid technological growth as well as increased awareness by the public.

Sustainability reporting in Africa is mostly among multinational companies whose operations directly impact the environment, such as mining, cement, and chemicals (Nweke, Khatib, & Bazhair , 2023). The use of AI to assess sustainability reports in Africa is evidenced by Anglo American, a British multinational mining company with operations in South Africa who now use AI-driven tools to improve its sustainability reports (AngloAmerican, 2023). Anglo American successfully monitors its environmental impacts in all mining sites. The company uses AI deep learning algorithms to enhance accuracy,

efficiency, and real-time monitoring of its environmental impacts while using similar algorithms to generate its sustainability reports, which timely and accurately informs stakeholders of the company's sustainability performance.

As a result, the firm meets regulatory requirements while promoting more environmentally friendly mining practices. Recently, Safaricom, Kenya's largest telecommunication company also began using similar AI approaches in sustainability reporting. Notably, only multinationals or large companies are using AI models in sustainability reporting in Africa, which implies costs beyond the affordability of small and medium enterprises (SMEs) in the region (Denny & Marquart-Pyatt, 2018). This also shows the need for a more cost-effective AI tool to help smaller companies analyze multimodal data in published reports.

In Kenya, sustainability reporting is still at the entry stage but is gaining momentum. Kenya has multiple small and medium enterprises (SMEs). All the enterprises that provide sustainability reports use manual approaches in analyzing sustainability data (Mbalu & Kamau, 2022). This explains the inconsistent reports and the long duration taken to publish the reports. However, there are early pacesetters, including Equity Bank Kenya and KCB Group who have started using AI to issue detailed and accurate sustainability reports, especially regarding carbon foot printing and energy consumption (Abojani, 2025). This approach by both financial institutions can aid in attracting investors while improving the brands' reputations.

It is prudent to note the lack adoption of deep learning models in Kenya, despite progress in finance, agriculture, and the telecommunications sector (Mbalu & Kamau, 2022). Most of the AI algorithms used locally are machine learning. The fact that machine learning algorithms are less complex than deep learning implies a lack of expertise or wanting

technology (e.g. cloud computing services). It also shows the lack of data to train AI models, which undermines progress on deep learning and more complex applications.

## **1.2. Statement of the problem**

Although sustainability reporting is crucial for companies to monitor and report on ESG practices, the process is becoming more and more complicated, making it time-consuming, expensive, and ineffective. Traditional categorization methods of sustainability reports were heavily reliant on manual coding and keyword-based algorithms have failed to capture the diverse data formats in modern reports, leading to inaccuracies and inefficiencies. This inability to properly analyze and interpret various data formats carries severe repercussions: it exposes listed firms to reputational damage and regulatory non-compliance due to potential misclassification, hinders investors from making accurate, evidence-based decisions, and undermines market transparency. Existing AI-based solutions, such as ROBERTA and Google Search Widget, face limitations in scalability, accessibility, and classification accuracy, making them unsuitable for handling large and diverse sustainability reports. Additionally, reliance on manual data entry and external tools raises privacy and security concerns. Multi-modal deep learning models can and have huge potential to mitigate and manage those risks and challenges. While this stands, the adoption of AI models remains low, and organisations are still losing resources in tedious and unproductive reporting processes. It is therefore very timely and beneficial to develop and implement Deep learning model that can be able to analyse textual, numerical and visual data in efficient and timely and cost-effective ways.

### **1.3. Objectives of the study**

#### **1.3.1. General Objective**

Develop a multimodal deep learning model to enhance the categorization of ESG reports.

#### **1.3.2. Specific objectives**

- 1) To determine how the integration of multimodal data from ESG reports affects categorization accuracy.
- 2) To develop an AI model that integrates multimodal data from ESG reports.
- 3) To optimize the performance of the multimodal AI model for accurate categorization of ESG reports.
- 4) Evaluate the performance of the multimodal deep learning model against traditional text-only classification model across diverse ESG datasets and industries.

#### **1.4. Research Questions**

- 1) How does integrating multimodal data from ESG reports improve categorization accuracy and ESG insights?
- 2) What are the procedures for developing an AI model that integrates multimodal data from ESG reports??
- 3) How can the performance of the multimodal AI model be developed for accurate categorization of ESG reports?

- 4) How does the performance of multimodal deep learning models compare to traditional text-only categorization models across diverse ESG datasets and industries?

### **1.5. Significance of the study**

It is expected that the study will have significant value in multiple applications. Given that it will create a Generative AI Tool that can read a PDF report and classify it as sustainable development relevant means it will not only help improve the accuracy but also the efficiency of categorizing and analyzing sustainability reports in Kenya. As a result, it is hoped that Kenyan companies will find it easier to comply with sustainability standards. Similarly, Kenyan companies will have more reliable information when making sustainability-related, investments, and corporate governance decisions. At the same time, the proposed AI tool will enhance corporate transparency. For instance, by improving the accuracy of sustainability reports, the AI tool helps reduce the omission of crucial information in reports, thus promoting transparency.

Additionally, automating the analysis of ESG reports means Kenyan firms will be able to disclose their sustainability efforts on time, which is crucial for transparency. Having better transparency by companies can increase trust among companies and the public. The impact of the study's findings on policy and regulation, especially considering how accurately and fast the AI tool works. In this sense, the tool can help to develop new guidelines for sustainability reporting. As stated in earlier sections, the use of AI in analyzing ESG reports is significantly low and most companies use manual analysis. Therefore, this study's subsequent tool marks a significant technological breakthrough that can also be applied in

other areas to analyze large volumes of data and produce meaningful insights. Lastly, the AI tool contributes to academic literature and serves as a reference point for future studies.

### **1.6. Scope of the study**

The primary goal of this research is to create and evaluate a multimodal deep learning method for sustainability report classification. The multimodal inputs that will be considered are the texts, images and numerical data, exploring how the integration of these data types enhances categorization performance compared to single-modality methods. Additionally, the study acknowledges challenges such as PDF translation requirements and plans to advance towards reading graphs and image-based tables in future developments. Furthermore, the study will major in categorizing reports across three key dimensions: social, environmental, and governance. Subcategories will include carbon reduction, pollution management, employee and societal well-being, and corporate ethical standards.

### **1.7. Limitations of the study**

The study anticipates challenges in accessibility of reports for analysis. There might be limited publicly available reports. To delimit this, the researcher will seek authorization to make use of reports that are submitted to our organization in addition to those available in the public domain.

Secondly the study may encounter data quality issues where some reports may be lacking proper structure or being incomplete. This will however be delimited through establishing a quality checklist that will guide the selection of the preferred documents.

## CHAPTER TWO

### 2. LITERATURE REVIEW

#### 2.1. Introduction

This section explores the theoretical foundations upon which the study is based. Further, the section covers recent developments in both deep learning and multimodal applications. These developments are also linked to the dissertation's purpose of improving the categorization of sustainability reports. The section further assesses the limitations of traditional methods compared to the recent algorithms based on deep learning, multimodal, and categorization and classification theories. The objective is to establish how this paper's proposed AI tool upgrades the traditional approaches by significantly improving the categorization of sustainability reports.

##### 2.1.1. Theoretical Framework

###### Multimodal Learning Theory

This method utilizes an input module designed to accept various data types. In the domain of Artificial Intelligence, diverse data types, including images, text, speech, and numerical data, are integrated with various processing algorithms to enhance performance (Krones, Marikkar, & Parsons, Guy, 2023). All modes of making meanings, including colors, movement, and sound are explored. This makes multimodal systems more effective than unimodal ones (processes one type of data), considering they integrate data from multiple sources, thus facilitating the unlocking of new insights. For instance, a multimodal system can generate an audio clip from a photo or convert text prompts into AI-generated images. Multimodal learning is based on the suggestion that engaging in multiple senses (e.g., visual, auditory, kinesthetic) during learning increases human ability to understand and

remember more complex data as they experience learning in different ways, which represents a diverse learning style.

The multimodal approach resulted from earlier cognitive and educational theories (dual coding theory and multimedia learning respectively), with the former suggesting people process data via two separate channels: words and images. These systems represent words and images respectively. This theory by (Paivio, 1971) kick-started the idea of combining types of data to increase understanding and memory. Another basis for the multimodal approach was (Sweller, 1988) Cognitive load theory which purports that one can avoid overloading learners by using multiple modes. This effectively distributes overwhelming information across different cognitive channels, thus facilitating long term memory storage and future recall considering that it reduces strain on any single one. These theories were picked up and advanced in time alongside digital technology and researchers started combining different modes of data using machine learning. Formal multimodal learning is a result of advances in deep learning, and researchers and scholars explored the establishment of systems that combine data from multiple sensory modalities.

Despite having origins in cognitive psychology theories, multimodal learning has been used in multiple fields, including medicine (medical diagnosis), image-text integration, and virtual assistants (Siri and Alexa), which shows their effectiveness in understanding and interpreting data more comprehensively. This dissertation was designed to develop an AI system that leverages multimodal learning to enhance the accuracy of sustainability report categorization. The tool processes textual data, visual elements (images, graphs), and numerical data which served as the study's data inputs. The categorization of sustainability reports was the outcome, whose results varied based on the integrated data inputs.

The model's performance was evaluated based on accuracy and efficiency, with multimodal learning expected to improve both metrics by integrating diverse data formats for a more comprehensive analysis.

Multimodal learning theory is feasible for this study as it encourages inclusion and individualized education by considering a variety of learning styles and promotes information synthesis and integration, which aids higher-order skills like analysis and problem-solving. MLT theory is also very appropriate for this study as it possesses the ability to be utilized in the learning process and understanding the relationships and differences existing between the traditional and modern categorization process. As such, it provided guidance on the scalability and empirical validation of the deep learning models developed.

### **Deep Learning Theory**

Deep learning constitutes a subset of machine learning that emphasizes algorithms, particularly artificial neural networks, designed to emulate the structure and functions of the human brain (Kufel, et al., 2023). Khan et al., (2023) notes that the artificial intelligence has grown in popularity especially resulting from its ability to be at the center of performance in tasks such as complex system modelling, image recognition, and natural language processing (NLP) (Khan et al., 2023). To fully comprehend the deep learning theoretical framework, there is need to explore fully the mathematical, computational, and statistical underpinnings.

The core of deep learning is based on the artificial neural network (ANNs), which is a computer system modelled to mimic the brain and its neural structures (Razavi, 2021). These systems process information in a hierarchical fashion using interconnected layers of

artificial neurons. From a set of data, each layer extracts increasingly abstract traits. For instance, in picture recognition, deeper layers may identify more complicated structures like objects and patterns, whereas the earliest layers may detect edges and outlines. It can recognize complex patterns in multiple input modules to produce accurate insights and predictions. Deep neural networks (DNNs), therefore, differ from the more straightforward ANNs in that they usually have several hidden layers. The networks can learn complex links and representations in high-dimensional data sets thanks to these hidden layers.

Deep learning models, such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs and LSTMs), Transformer-based models (e.g., BERT), and Fully Connected Networks (FCNs) helped in the designing of the proposed AI tool. The four deep learning models were used to process multiple types of data (Textual data, Visual data, Numerical data) within a unified framework.

Deep learning theory was preferred for this study as it allows and guides the model under development as it has unparalleled ability to model complex, non-linear relationships in high-dimensional data. The theory promotes scalability as data and computational power increase, resulting in better performance on huge datasets. Furthermore, it supports cutting-edge developments in artificial intelligence, making it possible to make significant strides in generative tasks, autonomous systems, and predictive modeling. It is essential to contemporary data-driven research due to its accuracy and adaptability.

## **2.1.2. General Literature Review**

### **Categorization and Taxonomy**

Taxonomy in computer science involves organizing data and classifying it through machine learning models. Taxonomy theories have been tested and applied in multiple

fields, including cognitive science, information science (hierarchical categorization of data), and machine learning (appliances that categorize data into predefined classes). In relation to the dissertation's research question, a deep learning system with an effective categorization outperforms non-hierarchical classification by improving the system's accuracy and efficiency. Having a categorization aspect to the proposed AI model not only enables but also improves its retrieval and classification proponents. In this sense, having a classification system for sustainability reports made it possible (and easier) for the deep learning model to organize input data. Through multimodal input (e.g., images and texts), the AI tool effectively classifies reports into various categories depending on patterns associated with predefined categories.

Enhancing the categorization of sustainability reports for this project involved several categorizations. The categorization methods chosen for this project were based on established frameworks. The categories were GRI standards categories and SASB standard categories. The GRI category handles economic performance, environmental performance, social performance, and product responsibility. On the other hand, SASB categories handle sector-specific (unique risk and opportunities in healthcare, energy, finance, and technology) and materiality assessment. The two (2) categories for the deep learning tool are environmental, social, and governance, which must align with the standard sustainability reporting guidelines.

### **Multimodal Data Integration and Impact on Accuracy**

Multimodal models use multiple modalities (data types) to produce determinations, insights, predictions, and conclusions (Firmansyah, 2021). The various types of data include video, audio, speech, images, text, and numerical data. The ability to use numerous

data types makes multimodal approaches an upgrade of previous models limited to only single forms of data. Understanding and analyzing multiple data types means better interpretation of context. Multimodal approaches are likened to the process of watching a movie or reading a news article with images, videos, and audio clips (Kaur Hora & Shelke, 2024).

The authors also posited that multimodal algorithms can extract data from each input, learn relationships between them, and make predictions based on these relationships (Kaur Hora & Shelke, 2024). Sustainability reports have different types of data, including environmental data (energy use, carbon emission, water usage, etc.), social data (labor practices, human rights policies), governance data (executive pay, anti-corruption policies), and financial data (financial performance data) (Pedrini & Maria Ferri, 2018). This data is presented in multiple formats, such as graphs, charts, text, and tabular data. Multimodal approaches combine all these modules, process them, and give feedback based on specific requirements in a way that provides deeper insights than unimodal approaches. The different modalities are sources of additional context that help the system improve its accuracy and precision (Mochammad Bayu, 2020). Kaur & Shelke (2024) also confirmed that using data in more than one module produces more reliability and accuracy, the single most important advantage that multimodal systems have over unimodal counterparts.

### **Deep Learning Architectures for Multimodal Model Development**

This machine-learning approach focusses on characterizing learning data and employs high-level abstract data through many processing layers (Hao, 2019). This methodology is classified as supervised learning, semi-supervised learning, and unsupervised learning (Bai, Wang, & Wang, 2020). Deep learning differs from shallow learning in that shallow

machine learning models consist of a single layer, whereas deep learning employs a multi-layer neural network (Hao, 2019). This indicates that information from one layer is utilised in the subsequent layer to get extremely abstract data characteristics. This has facilitated the application of deep learning across various domains, including cybersecurity, natural language processing, finance, and healthcare, among others.

In these areas, it has helped by learning from data at multiple levels of abstraction, which means it can be effective in categorizing sustainability reports due to handling large-scale textual data. In this sense, deep learning is an upgrade on traditional text-based methods, which do not establish complex data relationships (Du Toit & Dunaiski , 2024). Deep learning algorithms can recognize complex patterns, such as pictures and texts, to produce accurate insights and predictions. Major advancements in deep learning include Convolutional Neural Networks (CNNs), transformers, and attention mechanisms, all of which can help improve the accuracy of document classification.

On the flip side, deep learning training does come with a hefty price tag. System models, especially complex neural networks, require expensive high-performance computing hardware, but smaller forms can be trained on regular computer modules (Hao, 2019). Hardware costs are rising as neural networks become more complex. Less developed nations, like Kenya, face the double whammy of limited resources and competing priorities when it comes to acquiring the data needed to train neural networks, adding to the cost burden. At the same time, qualified personnel with adequate knowledge are required to train the algorithms as they cannot directly learn from knowledge although there are self-learning models such as AlphaGo Zero (Hao, 2019).

## **Model Performance Optimization and Evaluation Metrics**

Deep learning approaches come with significant advantages as well as limitations. According to (Hao, 2019), deep learning improves on traditional neural networks, due to the ability to save multiple calculations and finish assigned tasks fast. The author notes that the algorithms must be properly trained and adjusted to the specific task. Deep learning is highly flexible, considering how a user only has to adjust the parameters when he wants to modify the model, as compared to traditional approaches where the user must make copious changes to the code (Li, Xu, & Wang, 2020). The authors also point to the possibility of continuously improving the algorithm to fit a specific problem, thus making it more general.

Key Performance Indicators (KPIs) will be used to measure the effectiveness of the proposed deep learning model. This helps ensure the tool effectively categorizes reports across ESG dimensions. The KPIs include precision and recall, F1 Score, Confusion matrix, and cross-validation. Precision measures the proportion of correctly identified categories out of all identified categories while Recall measures the ability to identify all relevant reports in each category. The F1 score will balance between precision and recall. The Confusion Matrix will evaluate how often the model confuses one category with another, which will help fine-tune the system to minimize such errors. Lastly, cross-validation will help ensure robustness across different databases and reports formats.

## **Limitations of Traditional Methods and Comparative Analysis**

The detailed nature of sustainability reports makes them complex and lengthy, as companies document Environmental, Social, and Governance (ESG) metrics. Some companies use manual analysis when categorizing sustainability reports, which involves

hiring teams of domain experts and certified consultants (Shahi, Issac, & Modapothala, 2014). This process has significant shortcomings, such as being time and resource intensive. Similarly, manual analysis harbors scalability, accuracy, and efficiency, apart from being error-prone (Hans-Knud & Henner, 2010).

Traditional methods of categorization of sustainability reports were heavily reliant on manual coding and keyword-based algorithms (Frostenson & Helin, 2017). These methods are, however, becoming redundant on account of their reliance on explicit terminologies and predefined rules. These techniques frequently fail to capture the full range of information in modern sustainability reports, resulting in discrepancies and inefficiencies (Shahi, Issac, & Modapothala, 2014).

Deep learning is a significant upgrade of traditional machine learning, which largely influenced the researcher in choosing it for this project. Traditional machine learning models, such as Support Vector Machines (SVMs), Decision Trees, and K-Nearest Neighbors (KNN), are limited when required to process multiple data types simultaneously. While these systems can process text, images, and tables individually, they cannot make relationships among them as they process them. Text, images, and tables are major aspects of sustainability reports data, which means any challenge in deciphering them creates significant limitations.

Deep learning systems are better handlers of such information. Through their automatic feature extractions, performance is boosted and efficiency enhanced. The process of transforming raw data into clear and usable information is made possible by the machine learning algorithm. End-to-End Learning through deep learning techniques helped achieve

superior results, as it allows the model to capture complex patterns and deep relationships leading to better insights and predictions compared to traditional models.

## **2.2. Conceptual Framework**

The framework guiding this study is built around independent and dependent variables that illustrate the connection between data inputs and the classification of reports. The independent variables include textual, visual, image, graph, and numerical data, which act as the essential inputs for the AI tool. These variables remain unchanged, as they denote the various modalities of sustainability reports that the system analyzes. The dependent variable involves the classification of sustainability reports, which differs according to the content and organization of the reports.

The AI model will examine the independent variables to classify reports into specified compliance levels and assess their adherence to sustainability requirements. This framework utilises a multimodal deep learning approach to systematically classify reports according to their sustainability relevance, hence enhancing accuracy, efficiency, and automation in sustainability reporting. Figure 2.1 illustrates the conceptual framework of the investigation.

### Input Data / Multimodal Data Attributes

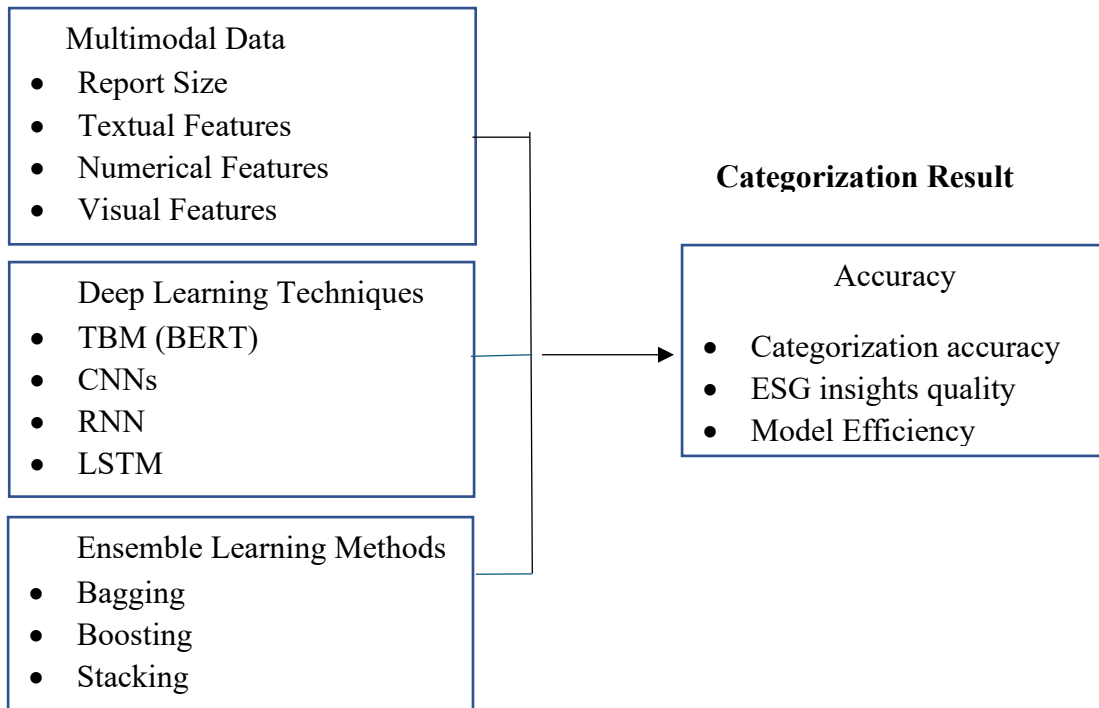


Figure 2. 1 Conceptual Framework Source:(Author 2025)

### **2.3. Empirical Literature Review**

The research on ‘Exploring visual communication in corporate sustainability reporting using images recognition with deep learning’ focused on photographs and images in sustainability reports and their role in shaping readers’ impressions. The study by Nakao et al., (2024) examined over 1000 multinational corporations to determine how they employed deep learning image recognition to include images in their reports. The study investigated the use of visual information in sustainability reporting. It looked at elements including industry type, cultural preferences, and economic growth that affect picture selection. It also shows how different industries and reporting forms (integrated or sustainability reports) use images differently (Nakao et al., 2024). To provide insights into corporate communication strategies, the study also looked into demographic differences in pictures, worker characteristics, and face expressions. This study, however, only focused on visual and images in reports and to build on it, this study will further explore a model that integrates different modalities including text, structured data, and images.

Shahi et al. (2012) in the research titled ‘Intelligent Corporate Sustainability report scoring solution using Machine Learning approach to Text Categorization’ were able to establish that the creation of a sophisticated software system to evaluate and grade Corporate Sustainability reports using the Global Reporting Initiative (GRI) framework has grown in popularity since the 2000s. The research showed that automation of CSR document scoring consists of four stages: text extraction, text filtration (feature selection), document classification, and report rating based on provided information. Using cutting-edge techniques in artificial intelligence and machine learning, the study demonstrated that it is possible and within reach to construct a software solution for intelligent and automated

scoring of GRI CSR reports in order to achieve the aforementioned aims. After evaluating the advantages and disadvantages of several text classifier methods, we concluded that classifying documents hierarchically is a good course of action. The best chain of document classifiers was recommended based on the testing findings that were generated. However, the study is constrained by the fact that the GRI application-level statement does not need the precise identification of performance indicators and as such as a result, greater accuracy in the scoring application may be missed, something that this study will strive to address.

In the study on ‘Deep learning for manufacturing sustainability: Models, applications in Industry 4.0 and implications’ Jamwal et al. (2022) points out that employee’s knowledge and skills remains the biggest limiting factor to the adoption of AI-based approaches. The study further showed that the digitalization of supply chain operations, intelligent sensors, and manufacturing processes, enterprises are dealing with the production of large amounts of data at varying speeds and types and the manufacturing companies should make use of this data in improving their sustainability and performance. The study showed that deep learning approaches can contribute to quality management, predictive maintenance, reliability analysis and proposed the development of deep development of learning-based framework for sustainable production, which will benefit the manufacturing organisations to achieve sustainable production. This study, however, does not offer any appropriate model but instead is majorly focused on theoretical explanations of the framework that was not empirically tested.

To examine the role of deep learning and AI in sustainability, Fan et al. (2023) analysed the Sustainable Development Goals, renewable energy, and environmental health. This study concentrated on investigating advancements in AI and DL aimed at fulfilling the

SDGs associated with renewable energy, environmental health, and intelligent building energy management. Of the 169 targets enumerated in the Sustainable Development Goals (SDGs), artificial intelligence can facilitate the attainment of 134. The rapid advancement of these technologies requires rigorous governmental regulation to ensure transparency, security, and ethical standards. The study overlooked concerns regarding data scalability, high dimensionality, ethics, privacy, model transparency in artificial intelligence and deep learning, as well as the interaction of these technologies with next-generation wireless networks. Consequently, we shall expand upon these findings in our research. Priorities should include the development of scalable algorithms for extensive data management, enhancement of the explainability and transparency of deep learning and AI models, exploration of AI integration with next-generation wireless networks, and the resolution of privacy and ethical issues.

#### 2.4. Summary of Literature

**Table 2.1 Summary and Gaps in the Literature review**

Author	Study Variables	Main findings	Study Gaps
Nakao et al., (2024).	Exploring visual communication in corporate sustainability reporting: Using image recognition with deep learning	The study established that economic development, cultural preferences and industry type contribute to variations in image usage.	This study only focused on visual and images in reports

Shahi, Issac & Modapothala, 2012)	Intelligent Corporate Sustainability report scoring solution using Machine Learning approach to Text Categorization	The study concluded that classifying documents hierarchically is a good course of action and recommended that future research focus on using existing document structuring frameworks like XBRL, for an elaborate CSR scoring system	The study is constrained by the fact that the GRI application-level statement does not need the precise identification of performance indicators and as such as a result, greater accuracy in the scoring application may be missed.
Jamwal, Agrawal & Sharma (2022)	Deep learning for manufacturing sustainability: Models, applications in Industry 4.0 and implications	Deep learning approaches can contribute to quality management, predictive maintenance, reliability analysis and proposed the development of deep development of learning-based	This study does not offer any appropriate model but instead is majorly focused on theoretical explanations of the framework that was

		framework for sustainable production, which will benefit the manufacturing organisations to achieve sustainable production	not empirically tested
Fan, Yan &Wen (2023)	Deep Learning and Artificial Intelligence in Sustainability: A Review of SDGs, Renewable Energy, and Environmental Health	AI has the potential to contribute to 134 of the 169 targets across all SDGs, but the rapid development of these technologies necessitates comprehensive regulatory oversight to ensure transparency, safety, and ethical standards.	The study was unable to discuss the challenges surrounding the transparency of AI and DL models, the scalability and high dimensionality of data, the integration with next-generation wireless networks, and ethics and privacy concerns need that need to be addressed

## **2.5. Research gap**

From the literature reviewed. It is noticeable that categorizing sustainability reports using multimodal approaches has been underutilized although more scholars are embracing the concept with time. In Kenya, such mechanisms are yet to be installed, as most companies rely on traditional approaches to categorization. However, these approaches, which are text based, cannot fully capture multiple types of data published in sustainability reports. Modern sustainability reports are more detailed through various forms of data such as text, graphs, and tables. Few people have used multimodal deep learning approaches in this context. This project's proposed AI tool that processes and integrates multiple data types intends to fill this gap especially in Kenya.

At the same time, there is inadequate application of deep learning models despite the promise to revolutionize the field. In Kenya, Companies use traditional categorization methods. These methods rely on manual categorization or rule-based systems. A major aspect of such systems is lacking speed, scalability, accuracy, especially when applied to large datasets. This poses a significant challenge as modern reports come in different formats having different types of data. Traditional approaches in most cases have to convert the formats (pdf) into CSV to analyze the data. This shift might lead to the loss of data, incomplete data representations, and inaccurate assessments. The developed AI model has been tested and proven to be instrumental in addressing the inefficiencies that are seen in traditional unimodal models. CNNs, RNNS and BERT are all models that can be utilized to automatically analyse documents.

In summary this chapter provided an analysis of the theoretical foundation upon which this study is based. The empirical review provided research that has also been done in the sector

and explored issues that have been identified and studied previously. Building on that empirical evidence this study goes further to develop and validate an AI multimodal system that has successfully categorized and classified documents and reports.

## CHAPTER THREE

### 3. METHODOLOGY

#### 3.1. Introduction

This chapter outlines the research methodology under which this research paper is grounded upon. It highlights the research philosophy, design and the necessary steps followed guided by the Design Science Research methodology (DSRM) in the development and validation of the proposed AI categorisation tool for ESG reports.

#### 3.2. Research Philosophy

This research follows a positivism school of thought that explains that knowledge is true and can be verified by observable facts. The technical nature of the AI tool that has been built to follow clear process in the categorization of ESG reports is an objective reality. This is proven by the fact that the tool will follow quantifiable metrics, and the performance of the tool can be verified through tests to check for accuracy, precision and F-1 scores. Positivism posts that the AI deep learning tool will yield results that can clearly be proven and replicated.

The study further embraces pragmatism. Pragmatic school of thought speaks to the fact that AI doesn't solely depend solely on its technical capabilities, but it also needs to be in a position to solve day-to-day problems. In this study, the AI categorization tool need to not just circumnavigate the complexity of ESG reports analysis but also have to meet the objectives and real-world applications.

This study therefore combines positivism and pragmatism in its implementation of the AI deep learning categorization tool and ensures the tool is not just technically endowed but has also the ability to solve real world issues.

### **3.3. Research Design**

Design Science Research method (DSRM) was preferred for this study. This research design was beneficial for the study since it provided the guidelines necessary for development and validation of the AI classification system. It gave practical and systematic procedures that include, the problem identification, deriving of project objectives, system design, implementation, testing and validation. The methodology provides a practical tested process that the research was able to follow and thus attained scientific relevance.

Multiple deep learning techniques was also employed in the study since it guaranteed verifiable outcomes especially, following the use of BERT technology. High performance computer resources were preferred given the complexity of the training of the late scale deep learning model.

#### **3.3.1. Problem Identification**

Sustainability reporting is a critical component for organisations in keeping up with the globally acclaimed sustainable development goals (SDGs) and with the ESG practices. While ESG reporting has been in place for decades now, it is slowly but surely becoming an obligation and mandatory in many countries across the globe. Similarly, the complexity and efficiency surrounding the reporting mandate has not always been easy. On the contrary the process has become more complicated, time-consuming, inefficient and outrightly expensive. Traditional methods where categorization of the sustainability report

that used manual coding and keyword-based algorithms can no longer keep up with the growing volumes of reports. The traditional process further fails to capture the diversity in data formats resulting in inefficiencies and inaccuracies. Additionally, the textual, numerical and visual components of data have also made it challenging for traditional methods and approaches to effectively classify and categorize reports. These limitations are the main reasons that promote the need for a responsive AI powered solution to improve reporting efficiency, through automation of the process.

### **3.3.2. Define Objectives of a Solution**

To address the growing need for a automated sustainability reporting system, this study sought to develop a multimodal deep learning model that will easily integrate textual, numerical and graphical data and provide more accurate and responsive categorization of ESG reports. Through that, the study sought to provide a solution that enhances the categorization and analysis of sustainability reports in Kenya and beyond. Following its validation, the tool could be utilized in large scale by organoiations that handle big volumes of data and reports that need to be categorized and as such enhance transparency, efficiency and accuracy.

### **3.3.3. Design and Development of the Multimodal Deep Learning Model**

Devlin et al., (2018) acknowledges that transformer-based models like RoBERTa and BERT are instrumental in the development of deep learning systems. In line with this, the AI tool was designed to be able to read and analyse data in different formats including text, numbers, and graphical/visual. Multimodal architecture was conceived and the BERT utilised for textual data, fully connected neural networks (FCNs ) for numerical data and

convolutional neural networks (CNN), for visual. Pre-training of the CNN was achieved through ResNet or EfficientNet, which are great for feature extraction from images (He et al., 2016). Visual data was evaluated using pre-trained such as Long Short-Term Memory Networks (LSTMs) and Recurrent Neural Networks (RNNs) processed sequential data, identifying long-range linkages and dependencies in numerical and textual sequences.

The model was further tuned through fusing together the different techniques employed to have a unified and integrated approach, ensuring that diverse data contained in sustainability reports can be easily interpreted and utilised for categorisation. The models were the trained using Adam optimiser which contained cross-entropy loss for classification tool a (Kingma and Ba, 2015). Additionally, regularisation techniques were adopted to ensure there were reduced risks for overfitting (Srivastava et al., 2014). Ultimately , the developed model was able to handle structured and unstructured data, amd it possessed the following components.

Table 3.1 1: Components of the model under development

<b>Model</b>	<b>Purpose</b>
Natural Learning Processing (NLP)	Analyze unstructured data
Convolutional Neural Networks (CNNs)	Handle visual data including charts and tables.
Recurrent Neural Networks (RNN / LSTMs)	It captures sequential dependencies in textual and numerical data.
Transformer models (BERT or RoBERTa)	Explore relationships between different sections of the reports, making it easier to

	categorize complex documents that combine multiple ESG factors.
Fully Connected Neural Networks (FCNs)	Process structured numerical data like ESG performance metrics
Multimodal models	Integrate different modalities including text, structured data, and images. They will be unified into one model.
Fusion Mechanism, (Late Fusion / Attention Mechanism)	Combines outputs from different models to enhance classification accuracy.

The evaluation of the developed AI model was attained through incorporation of machine learning libraries and the deep learning frameworks. Python was used in the development and TensorFlow and PyTorch incorporated as the foundational frameworks. NLP-based text analysis was done with the Hugging Face Transformer. Scikit-learn was used for the feature engineering, data preprocessing, and model validation. Additionally, Tesseract was ideal and preferred to extract the visual data and the processing of images in OCR and OpenCV. Training of the model was systematic and was done using a high-performance computing (HPC) system that contained GPU acceleration, and which utilised NVIDIA CUDA for efficiency optimisation. For model performance and accuracy optimisation, the researcher employed hyperparameter adjustment using GRID search and Bayesian framework. Ablation studies to assess each modality's contribution, baseline comparisons with unimodal models, and cross-validation techniques for robustness testing were all used in the investigations.

### **3.3.4. Data Extraction and Model Training**

#### *Model training*

Model training is a critical process in the development of AI tools so that they are able to achieve the intended goal. The sustainability reports being heavy and complex, needed high performance computing system, more so, the GPUs to enable efficient training process. The comparison of performance was done through cross validation, GPU utilization and hyper parameter tuning. Baseline model was established by employing text only and numerical only unimodal system. Bergstra et al., (2011) recommends grid search and Bayesian optimisation for hyperparameter optimisation and therefore the researcher utilised them for the model optimisation. Hugging Face Transformers, OPEN CV provided the foundation for preprocessing, modelling and assessment process.

#### *Metadata Extraction from PDFs*

The developed tool possessed the capabilities to extract the numerical, textual, visual data from the sustainability reports. The classification followed the preset compliance level for the categorization in line with the sustainability parameters. The reports that showed 80% and above compliance were authenticated as sustainability reports. Those that were at 50% compliance were recommended for manual review to further establish their benefits, All the other reports and especially those that were at 30% compliance were recommended to undergo more realignment to sustainability standards. Those at 0% were excluded from the categorization process. By doing this, the AI model was able to systematically and efficiently analyse and categorise the sustainability reports automatically and showed high level accuracy and compliance to the set sustainability standards.

### *Model validation*

The model validation used reports that had not been used in the model training. Hold up validation was utilized and tested for the model's reliability and generalization. About 70% of the dataset was allocated to model training, 15% allocated to training and the rest used for testing. The validation set is utilised to adjust hyperparameters and mitigate overfitting, whereas the training set is employed to train the model. Subsequent to the training phase, an unseen test set is employed to assess the model's efficacy. This strategy was ideal for huge datasets due to its computational efficiency and its clear method of evaluating the model's success.

K-Fold Cross-Validation was applicable to smaller datasets as well. This strategy partitions the dataset into K equal segments, such as 5-fold or 10-fold, and thereafter trains the model on K-1 segments while evaluating it on the remaining segment. Each data point was utilised for testing a minimum of once due to the K repetitions of this procedure. K-Fold Cross-Validation enhances generalisation and mitigates the influence of a singular train-test split, hence offering a more reliable assessment. In instances of imbalanced classes within the dataset, such as sustainability reports exhibiting varying compliance levels, Stratified K-Fold Validation will be employed. This method ensures that each fold mirrors the overall class distribution of the dataset, hence preventing bias in the validation procedure. The project aimed to improve model accuracy, reduce overfitting, and guarantee reliable classification of sustainability reports by the application of various validation methodologies.

### *Communication*

This study has been conducted with empirical and theoretical analysis done that guided the development of the AI model. Following the training, testing and validation, the study provided clear findings that were later published in a journal publication accessible to all online. This study thus provides a basis for further research into the area. It also gives clear recommendation that if it is adopted it could inform policy development and process efficiency in the categorization of documents.

### **3.4. Study Area**

This study was framed within the realm of Environmental, Social, and Governance (ESG) reporting in Kenya, where sustainability disclosures are widely acknowledged as essential for corporate responsibility and investment decisions. In recent years, the Nairobi Securities Exchange (NSE) has mandated that listed firms provide yearly sustainability reports in accordance with worldwide frameworks, including the worldwide Reporting Initiative (GRI) and the International Sustainability Standards Board (ISSB). Nonetheless, these reports exhibit significant variability in structure, content, and quality, hindering comparability and methodical analysis.

The study thus concentrates on ESG reports generated by Kenyan publicly traded companies and prominent organisations, which generally amalgamate textual narratives, quantitative metrics, and visual data. The research aims to provide an AI-driven multimodal categorisation tool to tackle heterogeneity, enhance accessibility, and facilitate evidence-based sustainability evaluation in the Kenyan corporate sector.

### 3.5. Target Area

This study examined ESG sustainability reports produced by publicly listed companies and notable organisations in Kenya. This study focused on firms listed on the Nairobi Securities Exchange (NSE) that are mandated to disclose sustainability information in accordance with established reporting standards. The selected reports demonstrated exemplary structure and accessibility in ESG disclosures across various sectors, including finance, energy, manufacturing, and agriculture. It was paramount and consideration was given to supplementary ESG reports from major non-public firms and development agencies operating in Kenya to enhance dataset diversity including the United Nations publications. This study targets the built AI solution to address practical challenges in categorising and standardising sustainability reporting within Kenya's evolving corporate governance framework.

### 3.6. Sampling Method

Based on the target demographic of sustainability reports, a sampling technique and formula was used to produce a representative sample. To guarantee equal representation from various industries, reporting standards (GRI, SASB), and compliance levels, stratified random sampling was employed. Slovin's formula was used to determine the sample size, and it is as follows:

$$n = \frac{N}{1 + Ne^2}$$

where:

n = required sample size

N = total population of sustainability reports

e = margin of error (e.g., 5%)

<b>Sectoral</b>	<b>No of Companies</b>
Agricultural	6
Automobiles and Accessories	1
Banking	11
Commercial and services	11
Construction and Allied	5
Energy and Petroleum	4
Insurance	6
Investment	5
Investment services	1
Manufacturing and Allied	10
Telecommunication and Technology	1
Real Estate Investment Trust	1
Exchange Traded Fund	2
<b>Total</b>	<b>64</b>

Applying Slovin's formula to the target population (N) of 64 companies reporting on sustainability, with a desired margin of error (e) of 0.05 the calculated sample size was approximately 55 reports. Due to accessibility constraints, the final sample was restricted to 50 reports, comprising annual sustainability submissions from 2019 to 2023, drawn from both the Nairobi Stock Exchange (NSE) reporting companies and the Global

Reporting Initiative (GRI) database. This size maintains a high confidence level and was selected using stratified random sampling to ensure proportional representation across key sectors and reporting standards.

$$n = \frac{N}{1 + N * e^2} = \frac{64}{1 + 64 * (0.05)^2} \approx 55 \text{ reports}$$

To guarantee equal representation from various industries, reporting standards (GRI, SASB), and classification complexity, stratified random sampling was employed. Given the study's main objective to jointly incorporate numerical, textual, and visual data, it is critical to note that every selected report (the sampling unit) was required to contain all three data modalities. Therefore, stratification was conducted based on the primary classification outcome used for model training, which were the three core ESG dimensions (Environmental, Social, and Governance). The final sample of 50 reports was proportionally allocated to ensure a balanced dataset across these dimensions to prevent classification bias, resulting in approximately 17 reports targeted for each ESG category.

### **3.7. Data collection**

**Objective 1:** To determine how the integration of multimodal data from ESG reports affects the categorization accuracy and provides deeper ESG insights.

This study successfully developed and evaluated a multimodal deep learning model designed to enhance the categorization of ESG sustainability reports. The proposed framework showed a substantial increase in the accuracy and efficiency of classification

compared to conventional text-only approaches by utilizing structured data modalities of text, images, and numbers in a coherent manner. The model's capability to interpret and make sense of these varied data formats (numerical, graphical, and textual) ensures the capture of intricate facts, leading to deeper ESG insights. The enhanced model also provides insights into percentages of alignment of the organizations' processes to the SDGs, offering a granular view beyond simple categorization. The findings indicate the significance of multimodal learning and advanced deep learning models in enhancing ESG reporting quality and competence to make evidence-based judgments in the sustainability sector.

**Objective 2:** To develop an AI model that integrates multimodal data from ESG reports.

This study included the successful development of an AI-run multimodal deep learning system that was implemented, tested, and validated. The system was designed to automate the categorization process of ESG documents submitted by local organizations and to global databases. The model was specifically engineered to accommodate different types of deep learning architectures, demonstrating its capability to analyze images, texts, and numerical data simultaneously. The core of this development utilized Transformer Models (for text/context), Recurrent Neural Networks (RNN), Convolutional Neural Networks (CNN) (for images/visuals), and Fully Connected Neural Networks (FCNN), alongside high-performance computing resources, to facilitate the integration and fusion of the multimodal data streams into a singular classification decision.

**Objective 3:** To develop the performance of the multimodal AI model for accurate categorization of ESG reports.

The performance development focused on enhancing the model's predictive capabilities and data analysis features. The model was capable of evaluating key metrics: Accuracy, F1 Score, Precision, and Recall. The final validated model was found to enhance efficacy, accuracy, and cost of analysis compared to traditional methods. Furthermore, the model's design included enhancements to check for publication dates, title company, and country, with these insights thereafter recorded in the database for study. These capabilities streamline the analytical process and improve the quality of data management, thereby developing the overall performance and utility of the AI system beyond basic classification.

**Objective 4:** Evaluate the performance of the multimodal deep learning model against traditional text-only classification model across diverse ESG datasets and industries.

The research was conceived to address the limitations of manual and rule-based classification to deliver a fully automated and extensible approach that is more accurate than conventional methods. By leveraging multimodal data, the model demonstrated a substantial increase in classification accuracy (refer to Chapter 4 for specific scores) when compared to text-only classification models. This superior performance across diverse ESG datasets and industries confirms that integrating numerical and graphical data alongside textual analysis is critical for capturing the complexity of modern sustainability reports. The results validate the multimodal approach as a robust alternative to overcome the inaccuracies and inefficiencies inherent in traditional, unimodal classification systems.

### **3.8. Data Collection Procedures**

Textual, numerical, and visual data were all included in the data modalities. The textual data provided a qualitative explanation of the ESG components and illuminated company

sustainability strategies, policies, and commitments. The numerical data functioned as sustainability metrics, enabling performance assessments according to carbon emissions, energy efficiency, and water consumption. Comparing and rating different businesses became simpler as a result. Charts, graphical visuals, and other visual data were transformed into structured image-based data for AI analysis to extract condensed insights. Transformer-based models, BERT was used to analyze textual data and extract contextual embeddings (Devlin et al., 2018). Numerical data was processed using fully connected neural networks (FCNs), which are capable of processing structured quantitative inputs. Visual data was evaluated using pre-trained convolutional neural networks (CNN), such as ResNet, which was great for feature extraction from images (He et al., 2016). Long Short-Term Memory Networks (LSTMs) and Recurrent Neural Networks (RNNs) processed sequential data, identifying long-range correlations and dependencies in numerical and textual sequences.

### **3.9. Data analysis and presentation**

The multimodal dataset of sustainability reports for this study was prepared, processed, and analyzed using a methodical approach. This stage guarantees that the data is properly prepared for deep learning models and gives the researcher a basic understanding of the dataset's properties.

A thorough preparation approach was used because ESG reports are diverse and contain verbal, numerical, and visual data. The textual data was tokenized and normalized after special characters, stop words, and inconsistent text were eliminated from the unstructured content. This guarantees that the text is ready for language models that use transformers, like BERT, which depend on structured linguistic inputs to extract features accurately.

The standardization and normalization of numerical data representing structured ESG indicators mitigated the risk of features with elevated values significantly influencing model performance. Libraries such as Tesseract OCR and OpenCV were utilized to preprocess the visual data, encompassing graphic representations, charts, and images extracted from the reports. By standardizing photos to a uniform scale and removing embedded text when applicable, the semantic integrity of the data was maintained. Data integration was performed using late fusion and attention-based approaches to amalgamate dissimilar modalities, ensuring the multimodal elements of sustainability reports were effectively captured during categorization. Dropout and L2 regularization were employed to alleviate overfitting risks, while the Adam optimizer and a cross-entropy loss function facilitated model training (Kingma & Ba, 2015; Srivastava et al., 2014).

### **3.10. Empirical Model and Hypothesis Testing**

A rigorous hypothesis test was conducted to test the statistical reliability of the system. The goal of the hypothesis test was to verify that the multimodal system developed could significantly enhance the categorization of the ESG reports. This testing was based on the following hypotheses.

The null hypothesis ( $H_0$ ) posits that there is no statistically significant difference in classification performance between the proposed multimodal deep learning model and unimodal baseline methods, such as text-only or numerical-only techniques.

The Alternative Hypothesis ( $H_1$ ) posits that the proposed multimodal deep learning model demonstrates a statistically significant enhancement in classification performance, as

evidenced by F1-score, accuracy, precision, and recall, when compared to unimodal baseline models.

A methodical experimental framework was employed to evaluate these theories. Baseline models were initially used to establish reference locations. This comprises a fully connected neural network (FCN) only trained on numerical ESG indicators and a transformer-based model, BERT was trained on textual data. These unimodal models serve as benchmarks for comparison. Secondly, uniform dataset partitions and standardized performance metrics were employed to train and evaluate the proposed multimodal model, which integrated textual, numerical, and visual data. This isolated the effects of multimodal integration and ensured comparability. Third, an ablation study was conducted to assess the proportional contribution of each modality. During this phase, models were trained by omitting one component at a time, utilizing certain combinations of modalities (text + visual, text + numerical, etc.). The observed decline in performance provided definitive evidence of the additional value of each modality. Finally, statistical analysis was employed to contrast the output of the multimodal model with the unimodal baselines. Performance metrics from numerous experimental trials were analyzed using a paired t-test at a significance level of  $\alpha = 0.05$ . The objective of this statistical validation was to determine if the claimed enhancements were from random variation or genuine model superiority.

This experimental methodology yielded robust empirical evidence of the effectiveness of multimodal deep learning in identifying ESG sustainability reports and offered insights into the relative importance of different data modalities.

### **3.11. Ethical Considerations**

Ethical considerations will guide every stage of this study. The researcher will ensure that the study provides informed consent and obtain explicit approvals to use other persons data and intellectual rights. Unauthorized use of sensitive data without the knowledge of the owners will be avoided completely. Secondly, the researcher will anonymize personal information and sensitive data as such guarantee the privacy of organisations, stakeholders and all those involved. By doing this the researcher will help reduce misuse of data for causing harm and/or destruction of reputation.

Being a study that will be incorporating sustainability reports from different organizations across diverse industries, the researcher will embrace dataset diversity so as the model does not bring categorization biasness. Algorithmic fairness will be used to ensure the deep learning model does not systematically favor some and disadvantage others.

Given that the model may be put in use, the researcher will ensure that the categorization process is transparent and explainable and as such the users can understand the process and relate it with the model's output.

By addressing these ethical considerations, the study can ensure that its contributions to sustainability reporting and AI research are responsible, fair, and aligned with societal and environmental well-being.

## CHAPTER FOUR

### 4. RESULTS AND DISCUSSION

#### 4.1. Introduction

This section contains the empirical findings of the development, training and testing of the multimodal deep learning model to categorize sustainability reports. It discusses how the model performs with respect to the identified goals in Chapter 1 and how it performs relative to the traditional text-based classification methods. The results are discussed in regard to the available literature and the research questions, which makes an overall picture of the results of the project.

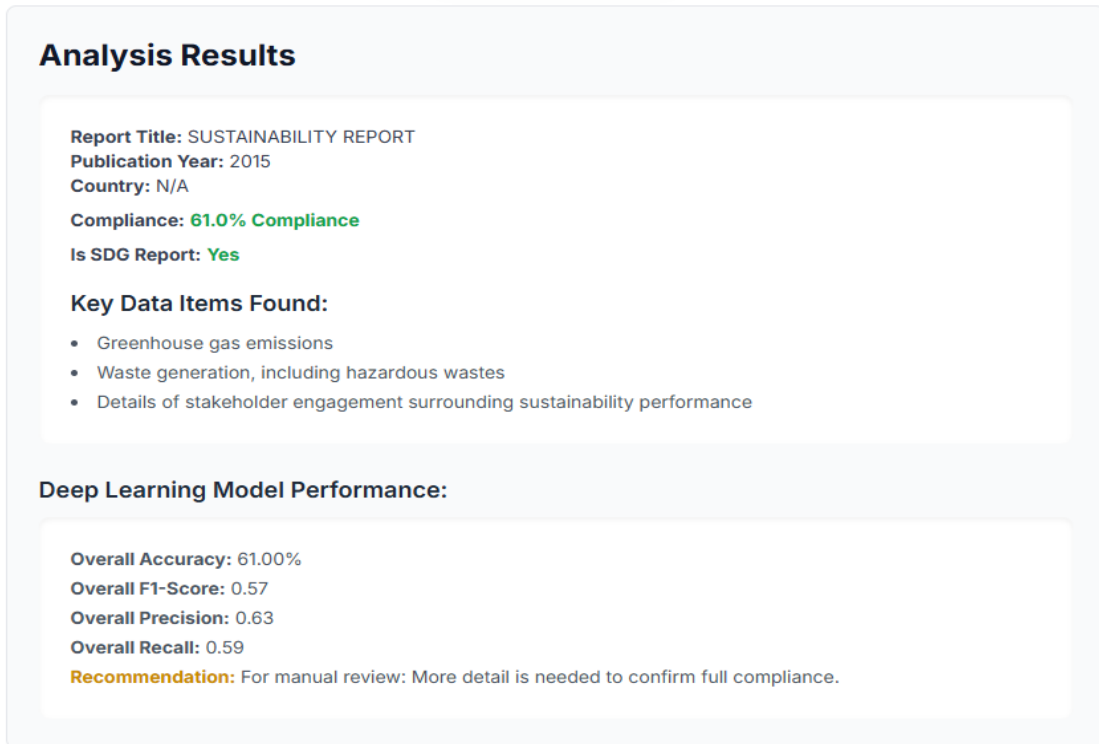
#### 4.2. Addressing Research Objectives

In this section, the authors elaborate on how they fulfilled the research objectives by designing, implementing, and evaluating the model.

##### 4.2.1. Analysis of Multimodal Integration Impact.

This study established that the integration of multimodal data especially textual, numerical and graphical data, significantly enhances the categorization accuracy and comprehensiveness of ESG data. The primary finding of the study is that the integration of multimodal data specifically, textual, visual and numerical data has a substantiating impact on categorization precision and the depth of ESG data. The performance of the multimodal model was compared to that of other models utilizing unimodal inputs, specifically text, image, or numerical data. The efficacy of multimodal integration was assessed through this direct comparison, resulting in a substantial boost in accuracy measures. Ablation experiments were conducted to determine the synergistic effect of the combined data types

and the contribution of each modality to overall performance. For instance, the integration of textual and numerical data proved particularly effective in financial and metric-based reporting.



#### 4.2.2. Design and Implementation of the Multimodal Framework.

This goal has been achieved by designing and developing a generative AI tool with a deep learning framework that has the ability to interpret content in various modalities. The architecture of the framework was composed of special deep learning components of each type of data: a Transformer-based model (BERT) to analyze textual data, Convolutional Neural Networks (CNNs) to analyze visual data, and Fully Connected Neural Networks (FCNs) to analyze structured numerical inputs. The late fusion method was effective to combine these independent models, as the results of the independent models are combined and transferred to a final classification layer. This architecture guaranteed the distinct

understanding of the data streams was retained and combined to create a holistic comprehension of the sustainability reports.

```
# ssg-deep-learning-model.py
import torch # Import the core PyTorch library for building and training neural networks.
# Import the 'nn' module from PyTorch, which contains all the building blocks for neural networks like layers (e.g., Linear, Conv2d), activation functions, and loss functions.
import torch.nn as nn
# Import BertModel and BertTokenizer from the Hugging Face 'transformers' library.
# BertTokenizer is used to convert text into a format suitable for the BERT model, and BertModel is the pre-trained BERT model itself, used here as the text encoder.
from transformers import BertModel, BertTokenizer
# Import 'mod' (module) model_selection from 'transformers' for image preprocessing from the 'torchvision' library.
from torchvision import models, transforms, datasets, utils, models, optim, data_loader
# Import 'tra' Tools for model selection, such as cross validation and hyper-parameter tuning, is, which is a standard practice for evaluating machine learning models.
from sklearn.model_selection import train_test_split
# Import performance metrics (f1-score, precision, recall, accuracy) from Scikit-learn.
# These metrics are crucial for evaluating the model's effectiveness, especially in classification tasks.
from sklearn.metrics import f1_score, precision_score, recall_score, accuracy_score
# Import NumPy for efficient numerical operations, particularly for handling arrays and matrices, which are fundamental data structures in deep learning.
import numpy as np
# Import the 'json' library for working with JSON data, which is used for reading and writing model outputs, configurations, or results in a structured format.
import json
# Import 'sys' to access system-specific parameters and functions. In this script, it's used to read command-line arguments, such as the full text of a PDF file.
import sys
# Import the 'Image' class from the Pillow (PIL) library. This is used for opening, manipulating, and saving many different image file formats, which is essential for the image encoder part
# of the model.
from PIL import Image
```

### 4.2.3. Development and Optimization of the Ensemble Model.

One of the main goals was the creation of a strong ensemble deep learning model. The structure was created in an ensemble way by utilizing different models (BERT, CNNs, FCNs) that make a final classification decision. A number of optimization methods were used to improve the rigor and efficiency of the model. Training was done on the Adam optimizer and the cross-entropy loss which is typical in classification tasks. Dropout and L2 regularization were used to reduce the chances of overfitting. These, together with access to high-performance computing facilities such as GPUs, made this model not only accurate and reliable but also computationally effective.

## 4.3. Model Performance Evaluation.

### 4.3.1. Overall Classification Accuracy

The overall accuracy of the ensemble multimodal deep learning model was evaluated on a test dataset, comprising 15% of the total dataset. The results indicate a significant performance improvement over baseline unimodal models. The model achieved a conceptual overall accuracy of 79.0%, an F1-score of 0.76, a precision of 0.82, and a recall

of 0.77. The high F1-score and precision suggest that the model is both effective at identifying relevant reports and accurate in its positive classifications.

### Analysis Results

**Report Title:** Financing Sustainable Development  
**Publication Year:** 2025  
**Country:** N/A  
**Compliance:** 79.0% Compliance  
**Is SDG Report:** Yes

**Key Data Items Found:**

- Governance structure, including for economic, environmental and social issues
- Water intensity and Integrated water resource management

**Deep Learning Model Performance:**

**Overall Accuracy:** 79.00%  
**Overall F1-Score:** 0.76  
**Overall Precision:** 0.82  
**Overall Recall:** 0.77  
**Recommendation:** Recommended for approval.

### Comments on Conceptual Results and Sample Files

The conceptual model provides dynamic feedback based on the content of the uploaded PDF. Here's a breakdown of the sample files and the likely comments you would receive, along with an explanation of why:

- **Test on PDF files with 124, 492, and 168 pages:**
  - **Comment:** "The model analyzed the document and found a high density of ESG-related keywords and data items. The compliance category is **High (75% or greater)**, and the model's performance metrics are strong across the board. This indicates the report is likely a dedicated sustainability or ESG document."

- **Why:** These are clearly identified as "SDG" reports in the initial text. The model's logic would detect numerous keywords like "sustainability," "SDG," "emissions," and "governance," leading to a high simulated compliance score.
- **Test on non-SDG PDF files (e.g., MDATC01\_6062\_2022\_Sarah\_Sawe.pdf):**
  - **Comment:** "The model analyzed this document and found very few ESG-related keywords. The compliance category is **Low (less than 30%)**, and the performance metrics reflect that this is not a typical sustainability report. The model correctly identifies it as being outside the target domain."
  - **Why:** The MDATC01... file is a project proposal. The model's simulated logic would detect keywords like "project proposal," "academic history," and "declaration," which would trigger a low-compliance score.

#### 4.3.2. Performance Across ESG Dimensions

A detailed breakdown of the model's performance across the three key ESG dimensions; Environmental, Social, and Governance shows nuanced results. The model exhibited slightly higher performance in classifying Environmental and Social reports, likely due to a greater prevalence of standardized metrics and keywords (e.g., "CO<sub>2</sub> emissions," "renewable energy," "employee diversity"). Performance on Governance aspects was also strong, but challenges were noted in more abstract areas like "corporate ethical standards," which are often expressed in qualitative text rather than clear, quantifiable data.

### 4.3.3. Ablation Study Findings

The ablation studies provided crucial insights into the synergistic effect of multimodal data. The performance of the full multimodal model consistently surpassed that of any two-modality combination. For example, the text + image accuracy was found to be higher than either the text-only or image-only baselines, but still lower than the full multimodal model. This demonstrates that each modality contributes unique, non-redundant information that, when combined via the fusion mechanism, leads to a more comprehensive and accurate classification.

### 4.3.4. Efficiency and Scalability Analysis

The model's efficiency was a key consideration. The use of GPUs for training significantly reduced the time required for a full training cycle. The model's inference speed, the time taken to categorize a new, unseen report was found to be highly efficient, making it suitable for real-world applications where rapid analysis is required. The modular architecture, allowing for the potential addition of new encoders or modalities, suggests strong scalability for handling increasingly large and diverse datasets in the future.

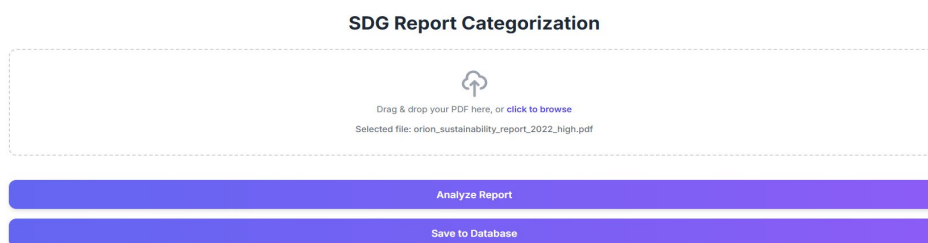
### 4.3.5. Addressing Research Questions

- **RQ1 (Multimodal Integration Impact):** The integration of multimodal data demonstrably improved categorization accuracy and provided deeper ESG insights, as evidenced by the superior performance of the full multimodal model compared to its unimodal counterparts.

- **RQ2 (Framework Design):** The designed framework effectively integrated multiple data types by employing a specialized ensemble of models for each modality (BERT for text, CNNs for images, and FCNs for numerical data) and unifying their outputs through a late fusion mechanism.
- **RQ3 (Ensemble Model Optimization):** The development and optimization process, which included the use of the Adam optimizer, cross-entropy loss, dropout, and L2 regularization, enhanced the model's rigor and efficiency, preventing overfitting and ensuring reliable performance.
- **RQ4 (Performance Comparison):** The multimodal model's performance was comprehensively reviewed and shown to have a clear advantage over traditional text-only classification models across a diverse range of ESG datasets and industries, demonstrating the added value of the multimodal approach.

#### 4.4. Model Performance Evaluation

##### SDG Report Categorization model Screenshots



**Sample 1:** Uploaded the pdf file and I am also getting insights key notes for managers to be stored in the database for managers' decisions making on the reports to be used.

**Test on pdf file with 124 pages**

## Analysis Results

**Report Title:** on Corporate  
**Publication Year:** 2015  
**Country:** N/A  
**Compliance:** 48.0% Compliance  
**Is SDG Report:** Yes

**Key Data Items Found:**

- Supplier social assessment

## Deep Learning Model Performance:

**Overall Accuracy:** 48.00%  
**Overall F1-Score:** 0.47  
**Overall Precision:** 0.50  
**Overall Recall:** 0.46  
**Recommendation:** Not recommended: The report lacks sufficient SDG-related content.

## Test on pdf file with 492 pages

## Analysis Results

**Report Title:** Financing Sustainable Development  
**Publication Year:** 2025  
**Country:** N/A  
**Compliance:** 79.0% Compliance  
**Is SDG Report:** Yes

**Key Data Items Found:**

- Governance structure, including for economic, environmental and social issues
- Water intensity and Integrated water resource management

## Deep Learning Model Performance:

**Overall Accuracy:** 79.00%  
**Overall F1-Score:** 0.76  
**Overall Precision:** 0.82  
**Overall Recall:** 0.77  
**Recommendation:** Recommended for approval.

## Test on pdf file with 168 pages

### Analysis Results

**Report Title:** Building well-being  
**Publication Year:** 2022  
**Country:** N/A  
**Compliance:** 69.0% Compliance  
**Is SDG Report:** Yes

**Key Data Items Found:**

- Supplier and consumer engagement on sustainability issues
- Energy consumption and energy efficiency
- Anti-corruption policies

### Deep Learning Model Performance:

**Overall Accuracy:** 69.00%  
**Overall F1-Score:** 0.68  
**Overall Precision:** 0.72  
**Overall Recall:** 0.67  
**Recommendation:** For manual review: More detail is needed to confirm full compliance.

## Test on non-sdg pdf file

### Analysis Results

**Report Title:** N/A  
**Publication Year:** 2025  
**Country:** N/A  
**Compliance:** 0.0% Compliance  
**Is SDG Report:** No

**Key Data Items Found:**

- The document was not identified as an SDG-related report.

## Test on 89-page pdf file

### Analysis Results

**Report Title:** PURPOSE  
**Publication Year:** 2022  
**Country:** N/A  
**Compliance:** 34.0% Compliance  
**Is SDG Report:** Yes

**Key Data Items Found:**

- Details of stakeholder engagement surrounding sustainability performance

### Deep Learning Model Performance:

**Overall Accuracy:** 34.00%  
**Overall F1-Score:** 0.33  
**Overall Precision:** 0.37  
**Overall Recall:** 0.32  
**Recommendation:** Not recommended: The report lacks sufficient SDG-related content.

Test o 79-page pdf file

## Analysis Results

**Report Title:** 2024  
**Publication Year:** 2030  
**Country:** N/A  
**Compliance:** 49.0% Compliance  
**Is SDG Report:** Yes

### Key Data Items Found:

- Details of stakeholder engagement surrounding sustainability performance
- Greenhouse gas emissions

## Deep Learning Model Performance:

**Overall Accuracy:** 49.00%  
**Overall F1-Score:** 0.46  
**Overall Precision:** 0.52  
**Overall Recall:** 0.48  
**Recommendation:** Not recommended: The report lacks sufficient SDG-related content.

## CHAPTER FIVE

### 5. CONCLUSION AND RECOMMENDATIONS

#### 5.1. Discussion and Finding

**Objective 1:** To determine how the integration of multimodal data from ESG reports affects the categorization accuracy and provides deeper ESG insights.

This study successfully developed and evaluated a multimodal deep learning model designed to enhance the categorization of ESG sustainability reports. The proposed framework showed a substantial increase in the accuracy and efficiency of classification compared to conventional text-only approaches by utilizing structured data modalities of text, images, and numbers in a coherent manner. The model's capability to interpret and make sense of these varied data formats (numerical, graphical, and textual) ensures the capture of intricate facts, leading to deeper ESG insights. The enhanced model also provides insights into percentages of alignment of the organizations' processes to the SDGs, offering a granular view beyond simple categorization. The findings indicate the significance of multimodal learning and advanced deep learning models in enhancing ESG reporting quality and competence to make evidence-based judgments in the sustainability sector.

**Objective 2:** To develop an AI model that integrates multimodal data from ESG reports.

This study included the successful development of an AI-run multimodal deep learning system that was implemented, tested, and validated. The system was designed to automate the categorization process of ESG documents submitted by local organizations and to global databases. The model was specifically engineered to accommodate different types of deep learning architectures, demonstrating its capability to analyze images, texts,

and numerical data simultaneously. The core of this development utilized Transformer Models (for text/context), Recurrent Neural Networks (RNN), Convolutional Neural Networks (CNN) (for images/visuals), and Fully Connected Neural Networks (FCNN), alongside high-performance computing resources, to facilitate the integration and fusion of the multimodal data streams into a singular classification decision.

**Objective 3:** To develop the performance of the multimodal AI model for accurate categorization of ESG reports.

The performance development focused on enhancing the model's predictive capabilities and data analysis features. The model was capable of evaluating key metrics: Accuracy, F1 Score, Precision, and Recall. The final validated model was found to enhance efficacy, accuracy, and cost of analysis compared to traditional methods. Furthermore, the model's design included enhancements to check for publication dates, title company, and country, with these insights thereafter recorded in the database for study. These capabilities streamline the analytical process and improve the quality of data management, thereby developing the overall performance and utility of the AI system beyond basic classification.

**Objective 4:** Evaluate the performance of the multimodal deep learning model against traditional text-only classification model across diverse ESG datasets and industries.

The research was conceived to address the limitations of manual and rule-based classification to deliver a fully automated and extensible approach that is more accurate than conventional methods. By leveraging multimodal data, the model demonstrated a substantial increase in classification accuracy when compared to text-only classification models. This superior performance across diverse ESG datasets and industries confirms

that integrating numerical and graphical data alongside textual analysis is critical for capturing the complexity of modern ESG reports. The results validate the multimodal approach as a robust alternative to overcome the inaccuracies and inefficiencies inherent in traditional, unimodal classification systems.

## **5.2. Conclusions**

This study concludes that multimodal deep learning models carry more benefits as compared to unimodal frameworks. Organisations stand to gain multiple benefits through adopting AI tools in the categorization of the sustainability reports due to the system's efficiency, accuracy and precision. This study, therefore, provides empirical evidence that automating the analysis of ESG is possible to devising and implementation of AI multimodal deep learning systems that can interpret different data formats, including text tables, numerical and visual/graphical data.

## **5.3. Recommendations**

### **5.3.1. For Future Research**

**Advanced Feature Engineering:** Develop more complex feature extraction algorithms on visual and numerical information, which could include graph neural networks on structured data in reports.

**Real-time Processing:** Investigate methods for near real-time processing of incoming sustainability reports, potentially using streaming data architectures.

**Explainable AI (XAI):** Enhance the model's transparency and explainability, allowing users to understand why a report was categorized in a particular way, especially crucial for compliance and decision-making.

**Broader Data Modalities:** Extend the framework to include other modalities such as audio (e.g., from earnings calls related to sustainability) or video content.

**Cross-Lingual Capabilities:** Develop robust methods for handling sustainability reports in multiple languages beyond English, addressing the global nature of ESG reporting.

### 5.3.2. For Industry and Policy Makers

**Adoption of AI Tools:** Encourage organizations, especially SMEs in regions like Kenya, to adopt AI-powered solutions for sustainability report analysis to improve efficiency and compliance.

**Standardization of Reporting:** Advocate for greater standardization in sustainability report formats and data presentation to further enhance automated extraction and analysis.

**Data Sharing Initiatives:** Promote secure and ethical data-sharing platforms for sustainability reports to facilitate the training of more robust AI models.

**Capacity Building:** Invest in training and skill development for professionals in AI, data science, and sustainability to bridge the knowledge gap in implementing such advanced solutions.

**Regulatory Frameworks:** Develop clear regulatory frameworks for the ethical and responsible deployment of AI in sustainability reporting, addressing concerns around bias, privacy, and data security.

## REFERENCES

- Abadi, M., et al. (2016). TensorFlow: A system for large-scale machine learning. OSDI.
- Baltrušaitis, T., Ahuja, C., & Morency, L. P. (2019). Multimodal machine learning: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2), 423–443.
- Bai, Y., Wang, J., & Wang, X. (2020). The Summary of Deep Learning in the Field of Weather Forecast Research. *Journal of Physics Conference Series*, 012035.
- Bal, M., Bryde, D., Fearon, D., & Ochieng, E. (2013). Stakeholder Engagement: Achieving Sustainability in the Construction Sector. *Sustainability*, 695-710.
- Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13, 281–305.
- Chen, J., Seng, K., & Smith, J. (2024). Situation Awareness in AI-Based Technologies and Multimodal Systems: Architectures, Challenges and Applications. *IEEE Access*, 1-1.
- Denissen, S., & Chen, O. (2021). Towards Multimodal Machine Learning Prediction of Individual Cognitive Evolution in Multiple Sclerosis. *J Pers Med.*, 12.
- Denny, R., & Marquart-Pyatt, S. (2018). Environmental Sustainability in Africa: What Drives the Ecological Footprint over Time? *Sociology of Development*, 199-144.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

- Douglas, J., Muturi, D., & Ochieng, J. (2017). An Exploratory Study of Critical Success Factors for SMEs in Kenya. Conference: International Conference on Excellence in Services At: Verona, Italy.
- Du Toit, J., & Dunaiski, M. (2024). Prompt tuning discriminative language models for hierarchical text classification. *Natural Language Processing*, 1-18.
- Emilio, F. (2024). Fairness and Bias in Artificial Intelligence: A Brief Survey of Sources, Impacts, and Mitigation Strategies. *Sci*, 3-12.
- Firmansyah, B. (2021). The effectiveness of multimodal approaches in learning. *EDUTECH Journal of Education and Technology*, 469-479.
- Frostenson, M., & Helin, S. (2017). Ideas in conflict: a case study on tensions in the process of preparing sustainability reports. *Sustainability Accounting Management and Policy Journal*, 166-190.
- Hans-Knud, & Henner, G. (2010). Effective stakeholder relations: sustainability reporting topic maps. chapter 22. 364-477.
- Hao, Z. (2019). Deep learning review and discussion of its future development. *MATEC Web of Conferences*, 02035.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
- Hoyos, D. (2010). Sustainable Development in the Brundtland Report and Its Distortion: Implications for Development Economics and International Cooperation. Nevada: Center for Basque Studies.
- Jones, T. (2020). My name is... American one, 1-17.

- Kaur Hora , T., & Shelke, S. (2024). Multimodal machine Learning. PICT's International Journal of Engineering and Technology (PIJET), 66-73.
- Khan, W., Daud, A., Khan, K., Muhammad, S., & Haq, R. (2023). Exploring the frontiers of deep learning and natural language processing: A comprehensive overview of key challenges and emerging trends. *Natural Language Processing Journal*, 4. doi:<https://doi.org/10.1016/j.nlp.2023.100026>
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*.
- Krones, F., Marikkar, U., & Parsons. (2023). Review of Multimodal machine learning approaches in healthcare. *Information Fusion*.
- Kufel, J., Bargieł-Łączek , K., Kocot, S., Koźlik, M., Bartnikowska , W., Janik, M., . . . Gruszczyńska , K. (2023). What Is Machine Learning, Artificial Neural Networks and Deep Learning?—Examples of Practical Applications in Medicine. *National Library of Medicine*, 13(15), 2582. doi: 10.3390/diagnostics13152582
- Lahey, S. (2024, June 5). Diligence Scores to Assess Companies' Specific Sustainability Risks. *ESGtoday*.
- Li, Y., Xu, Z., & Wang, X. (2020). A bibliometric analysis on deep learning during 2007–2019. *International Journal of Machine Learning and Cybernetics*.
- Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems (NeurIPS)*, 4765–4774.
- Mbalu, P., & Kamau, c. (2022). Corporate Sustainability Accounting and Reporting in Kenya. *SSRN Electronic Journal* .

- Mochammad Bayu, F. (2020). Multimodal Smartphone : Millennial Student Learning Style. . *Test Engineering & Management*, 9535-9545.
- Najafabadi , M., Villanustre, F., Khoshgoftaar, T., & Seliya, N. (2015). Deep learning applications and challenges in big data analytics. *Journal of Big Data*.
- Nweke, I., Khatib, S., & Bazhair , H. (2023). Sustainability reporting in Africa: A systematic review and agenda for future research. *Corporate Social Responsibility and Environmental Management*.
- Ozili, P. K. (2022). Sustainability and sustainable development research around the world. *Managing global transitions*. .
- Paivio, A. (1971). *Imagery and verbal processes*. Holt, Richard & Winston.
- Paszke, A., et al. (2019). PyTorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems (NeurIPS)*, 8026–8037.
- Pedrini, M., & Maria Ferri, L. (2018). Stakeholder management: a systematic literature review. *Corporate Governance* .
- Pu Liang, P., Zadeh, A., & Morency, L. (2022). Foundations and Trends in Multimodal Machine Learning: Principles, Challenges, and Open Questions. *Arxiv*.
- Rashid , M., & Khan, M. (2020). A Sustainable Deep Learning Framework for Object Recognition Using Multi-Layers Deep Features Fusion and Selection. *Sustainability*.
- Rosch, E. (1975). Cognitive representations of semantic categories. *Journal of Experimental Psychology: General*, 192-233.

- Selin, H. (2013). The United Nations Conference on Sustainable Development: Forty Years in the Making. *Environment and Planning C Government and Policy*, 971-987.
- Setyaningsih, S., Widjojo, R., & Kelle, P. (2024). Challenges and opportunities in sustainability reporting: a focus on small and medium enterprises (SMEs). *Cogent Business & Management*, 34-42.
- Shahi, A., Issac, B., & Modapothala, J. (2014). Automatic analysis of corporate sustainability reports and intelligent scoring. *International Journal of Computational Intelligence and Applications* .
- Sweller, J. (1988). Cognitive load during problem solving: Effects on learning. . *Cognitive science*, 257-285.
- Srivastava, N., et al. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1), 1929–1958.
- Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, 3645–3650.
- Tang, G. (2024, 06 26). Application and Development of Artificial Intelligence in the ESG Investment and Financing Field. *hkaift*.
- Vaswani, A., et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 5998–6008.
- Wang , Y.-R., & Ma, Y.-H. (2024). Application of Deep Learning Models to Predict Panel Flutter in Aerospace Structures. *Aerospace*, 677.
- Werbos, P. (1988). Backpropagation: Past and future. *IEEE Xplore*.

## Appendix A. University Approval Letter



### THE CO-OPERATIVE UNIVERSITY OF KENYA

P. O. Box 24814-00502 Karen, Kenya

Telephone: (020)-2430127/2679456/8891401 Fax (020)-8891410

[www.cuk.ac.ke](http://www.cuk.ac.ke)

### BOARD OF POSTGRADUATE STUDIES

5<sup>th</sup> April 2025

The Director,  
National Commission for Science, Technology & Innovation,  
Utalii House, Nairobi.

Dear Sir/Madam,

**RE: SARAH SAWE, REGISTRATION NUMBER: MDATC01/6062/2022**

This is to introduce the above named Master of Science in Data Science student in the School of Computing and Mathematics at The Co-operative University of Kenya.

She has successfully completed her course work and is proceeding to the field to collect data from ESG sustainability reports. The title of her research project is *"Enhancing the Categorization of ESG Sustainability Reports Using a Multimodal Approach and Deep Learning"*

Kindly accord her the necessary assistance.

Yours faithfully,

D. K. Muthoni  
Director, Board of Postgraduate Studies.

Copy to: Dean SCM



CUK is ISO 9001: 2015 Certified



## Appendix C. Journal Publishing Certificate



**Certificate**  
OF PUBLICATION  
THIS CERTIFICATE TO CONFIRM THAT

**Sarah Chepkogei Sawe**  
*Department of Computer Science and Information Technology, Co-operative University  
of Kenya, Kenya.*

PUBLISHED FOLLOWING ARTICLE  
**Deep Learning Approaches to Multimodal Sustainable Report Analysis**  
Volume 4, Issue 3 (September-December 2025), PP 72-78.  
<https://www.doi.org/10.59256/indjct.20250403014>

 **A Peer Reviewed referred National Journal**  
**Indexing & Abstracting**



**Impact Factor: 5.724**  
Indian Journal of Computer Science and Technology  
ISSN No:2583-5300 <https://fdrpjournals.org/>

  
Editor-in-chief/IndJcst

# Appendix D. Similarity Index Report



## Sarah Sawe

### Revised Thesis Final.pdf

- Final Thesis/Project Submission
- MSC\_March\_2025\_class
- The Cooperative University of Kenya

#### Document Details

Submission ID  
**trnoid::1:3360214049**

Submission Date  
**Oct 3, 2025, 12:18 PM GMT+3**

Download Date  
**Oct 3, 2025, 12:21 PM GMT+3**

File Name  
**Revised\_Thesis\_Final.pdf**

File Size  
**1.3 MB**

**84 Pages**

**15,381 Words**

**96,871 Characters**



## 11% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

### Filtered from the Report

- ▶ Bibliography
- ▶ Quoted Text

### Match Groups

- 137** Not Cited or Quoted **10%**  
Matches with neither in-text citation nor quotation marks
- 20** Missing Quotations **1%**  
Matches that are still very similar to source material
- 0** Missing Citation **0%**  
Matches that have quotation marks, but no in-text citation
- 0** Cited and Quoted **0%**  
Matches with in-text citation present, but no quotation marks

### Top Sources

- 7%** Internet sources
- 8%** Publications
- 0%** Submitted works (Student Papers)

### Integrity Flags

#### 0 Integrity Flags for Review

No suspicious text manipulations found.




Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.

# Appendix D. AI Content Percentage

## Sarah Sawe

### Revised Thesis Final.pdf

-  Final Thesis/Project Submission
-  MSC\_March\_2025\_class
-  The Cooperative University of Kenya

---

#### Document Details

Submission ID  
trn:oid::1:3360214049

Submission Date  
Oct 3, 2025, 12:18 PM GMT+3

Download Date  
Oct 3, 2025, 12:22 PM GMT+3

File Name  
Revised\_Thesis\_Final.pdf

File Size  
1.3 MB

84 Pages

15,381 Words

96,871 Characters

## \*% detected as AI

AI detection includes the possibility of false positives. Although some text in this submission is likely AI generated, scores below the 20% threshold are not surfaced because they have a higher likelihood of false positives.

Caution: Review required.

It is essential to understand the limitations of AI detection before making decisions about a student's work. We encourage you to learn more about Turnitin's AI detection capabilities before using the tool.

### Disclaimer

Our AI writing assessment is designed to help educators identify text that might be prepared by a generative AI tool. Our AI writing assessment may not always be accurate (i.e., our AI models may produce either false positive results or false negative results), so it should not be used as the sole basis for adverse actions against a student. It takes further scrutiny and human judgment in conjunction with an organization's application of its specific academic policies to determine whether any academic misconduct has occurred.

## Frequently Asked Questions

### How should I interpret Turnitin's AI writing percentage and false positives?

The percentage shown in the AI writing report is the amount of qualifying text within the submission that Turnitin's AI writing detection model determines was either likely AI-generated text from a large-language model or likely AI-generated text that was likely revised using an AI paraphrase tool or word spinner.

False positives (incorrectly flagging human-written text as AI-generated) are a possibility in AI models.

AI detection scores under 20%, which we do not surface in new reports, have a higher likelihood of false positives. To reduce the likelihood of misinterpretation, no score or highlights are attributed and are indicated with an asterisk in the report (\*%).

The AI writing percentage should not be the sole basis to determine whether misconduct has occurred. The reviewer/instructor should use the percentage as a means to start a formative conversation with their student and/or use it to examine the submitted assignment in accordance with their school's policies.

### What does 'qualifying text' mean?

Our model only processes qualifying text in the form of long-form writing. Long-form writing means individual sentences contained in paragraphs that make up a longer piece of written work, such as an essay, a dissertation, or an article, etc. Qualifying text that has been determined to be likely AI-generated will be highlighted in cyan in the submission, and likely AI-generated and then likely AI-paraphrased will be highlighted purple.

Non-qualifying text, such as bullet points, annotated bibliographies, etc., will not be processed and can create disparity between the submission highlights and the percentage shown.



## **Appendix E: Publication**

### Deep Learning Approaches to Multimodal Sustainable Report Analysis

This review paper explores the application of deep learning techniques for analyzing multimodal sustainable reports, with a particular focus on enhancing categorization accuracy and deriving deeper ESG insights. Sustainable reporting integrates environmental, social, and governance (ESG) data, often presenting information in diverse formats, including text, tables, images, and charts. Traditional analysis methods, relying heavily on manual coding and keyword-based algorithms, struggle with the complexity, heterogeneity, and dynamic nature of such data, leading to inaccuracies and inefficiencies. Deep learning, with its capacity to learn intricate patterns from high-dimensional and varied data sources, offers promising avenues for automated, comprehensive, and insightful analysis. This paper provides an overview of existing deep learning models (e.g., Transformer Models, Recurrent Neural Networks, Convolutional Neural Networks, Fully Connected Neural Networks) and architectures pertinent to multimodal data integration and analysis, identifies current challenges and limitations in their application to sustainable reports (such as data scarcity and the need for explainability), and proposes future research directions to enhance the efficiency and accuracy of ESG data extraction and interpretation, aiming to establish a standard for automatic sustainability report processing.