

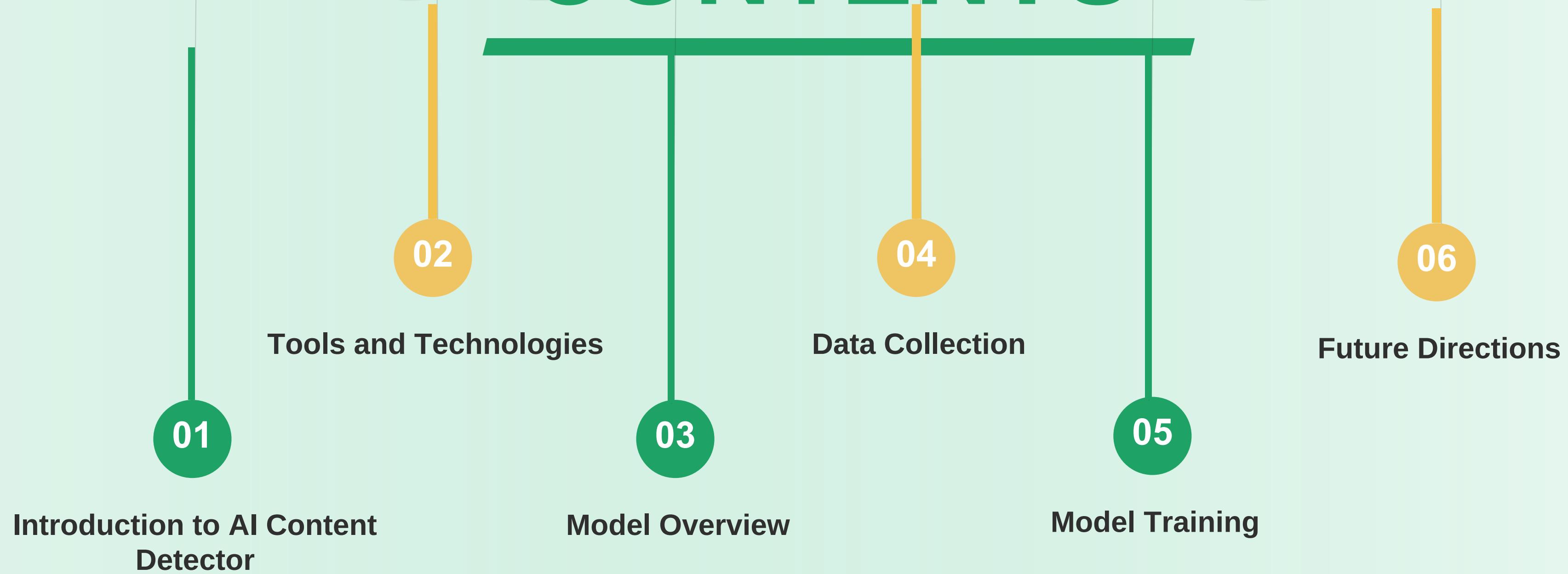
# AI Content Detector



Khan Talha  
Khan Imran  
Jaideep Yadav  
Khan Ejaj

In the Guidance of Er. ShridnidhiGIndi

# CONTENTS



# 01

## Introduction to AI Content Detector



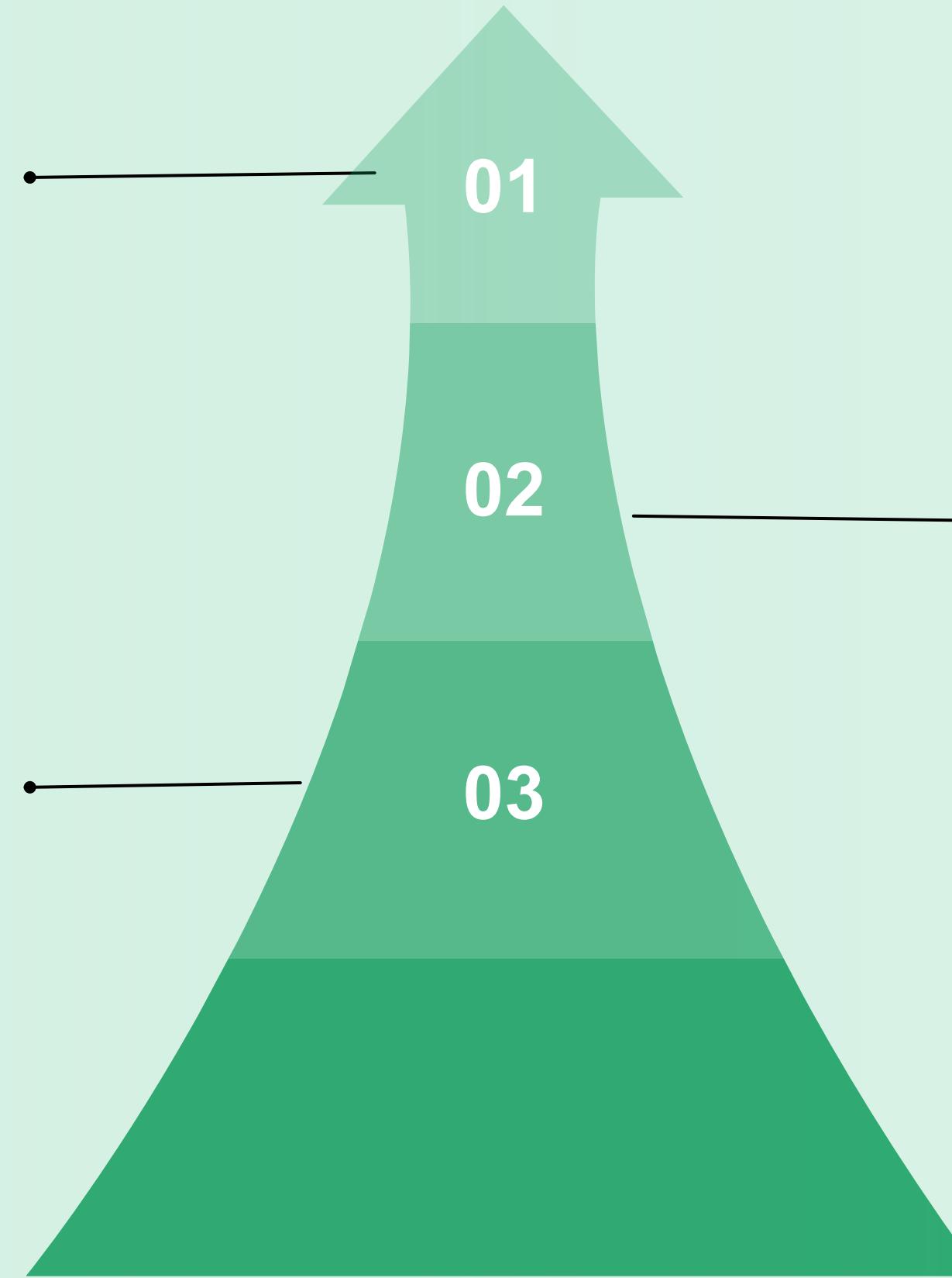
# Purpose and Scope

## Objective

The objective is to create a machine learning model that accurately differentiates between human and AI-generated content, improving content reliability and authenticity.

## Importance in various fields

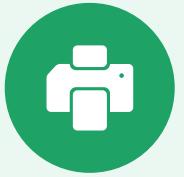
This model is crucial in fields like education, where integrity and originality are paramount, and journalism, where verifying sources is essential.



## Target users

The primary users include educators, journalists, and content moderators who require tools to assess content authenticity and mitigate misinformation.

# Applications



## Education sector

In education, the model can be utilized to detect plagiarism and ensure that students submit original work, maintaining academic integrity.



## Journalism practices

Journalists can use the model to authenticate sources and enhance credibility by distinguishing between human reporting and AI-generated stories.



## Content moderation

Content moderators can employ this tool to filter out misleading or automatically generated posts, enhancing the quality of online discourse.

# LITERATURE SURVEY

PAPER	YEAR	PURPOSE	TECHNIQUE	RESULT	PARAMETERS	FINDINGS	DRAWBACKS
Effectiveness of Free Software for Detecting AI-Generated Writing	2023	To explore the effectiveness of freely available AI detection software in identifying AI-generated content in student assignments and to assess its implications for teacher evaluation.	<ul style="list-style-type: none"><li>Manual analysis of student writing assignments.</li><li>Comparison of results from five different free AI writing detectors.</li></ul>	<ul style="list-style-type: none"><li>The findings revealed limitations in the effectiveness of free AI detection tools, indicating that they may not reliably identify AI-generated content.</li></ul>	<ul style="list-style-type: none"><li>Participants: Japanese university students writing in English as a Foreign Language (EFL).</li><li>Tools: Five different free AI writing detection software.</li></ul>	<ul style="list-style-type: none"><li>Highlighted the challenges educators face in detecting AI-generated writing.</li><li>It emphasized the need for cautious use of detection tools in educational settings, as they may not provide definitive conclusions regarding student honesty.</li></ul>	<ul style="list-style-type: none"><li>The free tools tested showed limitations in accuracy and reliability.</li><li>The evolving sophistication of AI-generated writing makes it increasingly difficult for detection algorithms to keep pace.</li></ul>
DeepFakeNet: A Deep Learning Approach	2019	<ul style="list-style-type: none"><li>To survey existing literature on deepfake detection.</li><li>Validate selected studies on deepfake detection.</li><li>Analyze empirical evidence in deepfake detection.</li></ul>	<ul style="list-style-type: none"><li>Various detection techniques (e.g., machine learning, deep learning).</li><li>Quality assessment criteria.</li><li>Data extraction and synthesis.</li></ul>	<ul style="list-style-type: none"><li>Performance evaluation of detection methods.</li><li>Finalized 91 research articles and 21 reviews.</li><li>Accumulated 112 studies.</li></ul>	<ul style="list-style-type: none"><li>Measurement metrics (e.g., accuracy, precision, recall).</li><li>Criteria for inclusion.</li><li>Publication period, methods, datasets</li></ul>	<ul style="list-style-type: none"><li>Identified trends in deepfake detection research and tools.</li><li>Comprehensive overview of existing research.</li><li>Emergence of deepfake research since 2018.</li></ul>	<ul style="list-style-type: none"><li>Limited datasets and variability in performance across methods.</li><li>Potential bias in selected studies.</li><li>Rapidly evolving technology outpacing detection methods.</li></ul>

# LITERATURE SURVEY

PAPER	YEAR	PURPOSE	TECHNIQUE	RESULT	PARAMETERS	FINDINGS	DRAWBACKS
Watermarking Techniques for AI-generated Images	2023	The study aims to investigate the robustness of watermarking methods against common post-processing techniques that can be used to evade detection of AI-generated images.	<ul style="list-style-type: none"><li>Watermarking methods: HiDDeN and UDH.</li><li>Post-processing methods: Gaussian blur, Brightness/Contrast adjustments, JPEG compression, and Gaussian noise.</li><li>The study also introduces an adversarial post-processing method called WEvade-W-II.</li></ul>	<ul style="list-style-type: none"><li>The study finds that existing post-processing methods can effectively evade watermark detection, and the proposed WEvade-W-II outperforms these traditional methods in terms of evasion rate while maintaining image quality.</li></ul>	<ul style="list-style-type: none"><li>For Gaussian blur: Kernel size <math>s=5</math> and varying standard deviation <math>\sigma</math>.</li><li>For Brightness/Contrast: Parameters <math>a</math> (scaling factor) and <math>b=0.2</math> (offset).</li><li>Watermark lengths: 30-bit for HiDDeN and 256-bit for UDH.</li><li>Detection thresholds <math>\tau</math> for evaluating the performance of detectors.</li></ul>	<ul style="list-style-type: none"><li>The evasion rate of post-processed watermarked images is significant, indicating that common image manipulations can undermine watermark detection.</li><li>The double-tail detector shows higher FPR compared to the single-tail detector, particularly at lower thresholds.</li></ul>	<ul style="list-style-type: none"><li>Theoretical FPRs do not exactly match empirical results due to watermark selection randomness.</li><li>Watermarking methods may be vulnerable to sophisticated attacks.</li><li>Specific parameter settings limit generalizability across datasets and applications.</li></ul>
Evaluating the Efficacy of AI Content Detection Tools in Differentiating Between Human and AI-Generated Text	2023	<ul style="list-style-type: none"><li>To assess the ability of AI detection tools to accurately identify AI-generated text from human-written content.</li><li>To explore the challenges AI-generated content poses to academic integrity.</li></ul>	<ul style="list-style-type: none"><li>Various detection techniques (e.g., machine learning, deep learning) used to differentiate between human-written and AI-generated text.</li><li>AI detection tools identified GPT 3.5-generated content more accurately than GPT 4 content.</li><li>Tools showed inconsistencies with human-written control responses, producing false positives.</li></ul>	<ul style="list-style-type: none"><li>Measurement metrics (e.g., accuracy, precision, recall).</li><li>Criteria for inclusion.</li><li>Publication period, methods, datasets</li></ul>	<ul style="list-style-type: none"><li>Identified trends in deepfake detection research and tools.</li><li>Comprehensive overview of existing research.</li><li>Emergence of deepfake research since 2018.</li></ul>	<ul style="list-style-type: none"><li>Potential bias in selected studies.</li><li>Rapidly evolving technology outpacing detection methods.</li></ul>	

# 02

## Tools and Technologies



# Tools and Technology

## Software Used

- **Python:** simple syntax, powerful libraries for efficient machine learning development.
- **Jupyter Notebook:** easy prototyping, interactive data visualization, and seamless sharing.
- **VS Code:** robust, extensible environment for efficient development, debugging, and collaboration.

## Libraries:

- **Transformers:** A library by Hugging Face for state-of-the-art natural language processing.
- **Torch:** The core library for PyTorch, used for machine learning and deep learning models.
- **Gunicorn:** A Python WSGI HTTP server for deploying web applications, useful for serving your models in production environments.

## Models:

- **RoBERTa:** Robustly Optimized BERT Pretraining Approach is a variant of BERT that improves performance by training longer, with larger batches, and removing the Next Sentence Prediction task.

# 03

## Model Overview



# Model Choice

- **Introduction to RoBERTa**

RoBERTa, a variant of BERT, is designed for enhanced performance in tasks involving natural language understanding, enabling nuanced comprehension of context.

- **Why RoBERTa?**

Its ability to analyze text at a deeper level makes RoBERTa especially suitable for distinguishing between the subtleties of human versus AI-generated content.

For Example:

"The bank is situated on the \_\_\_\_\_ of the river.



# Why Use Roberta?

- **Linguistic Features**

AI Content: Climate change is a critical issue. It affects weather patterns. Solutions include renewable energy. Governments must act now.

Simplicity: Shorter sentences, basic vocabulary (e.g., "talk to people far away," "play games"), and conversational tone.

Subjectivity: Personal advice ("listen to our parents") and emotional cues ("enjoy nature").

Subtle Errors: Less precise phrasing (e.g., "waste time" vs. "raise concerns about productivity loss").

- **Structural Patterns**

Human Essay:

Narrative Flow: Casual structure with abrupt topic shifts (e.g., jumping from emergencies to homework).

Repetition: Reiterates basic ideas ("use mobile phones wisely") without elaboration.

AI Essay:

Logical Progression: Each paragraph introduces a distinct theme (functionality → apps → sensors → societal debates).

Balance: Neutral, comprehensive coverage of pros/cons (e.g., "versatility" vs. "privacy concerns").

# 04

## Data Collection



# Types of Data



01

## Human-generated content

Human-generated content is sourced from blogs, articles, and academic papers, representing authentic expressions and thought processes of individuals.

02

## AI-generated content

AI-generated content is collected from automated writing tools and other generative models, allowing for a balanced dataset to train and validate the model.

# Challenges and Solutions



## Common challenges

Common challenges include dealing with imbalanced datasets and overfitting, which may hinder the generalizability of the model to unseen data.



## Mitigation strategies

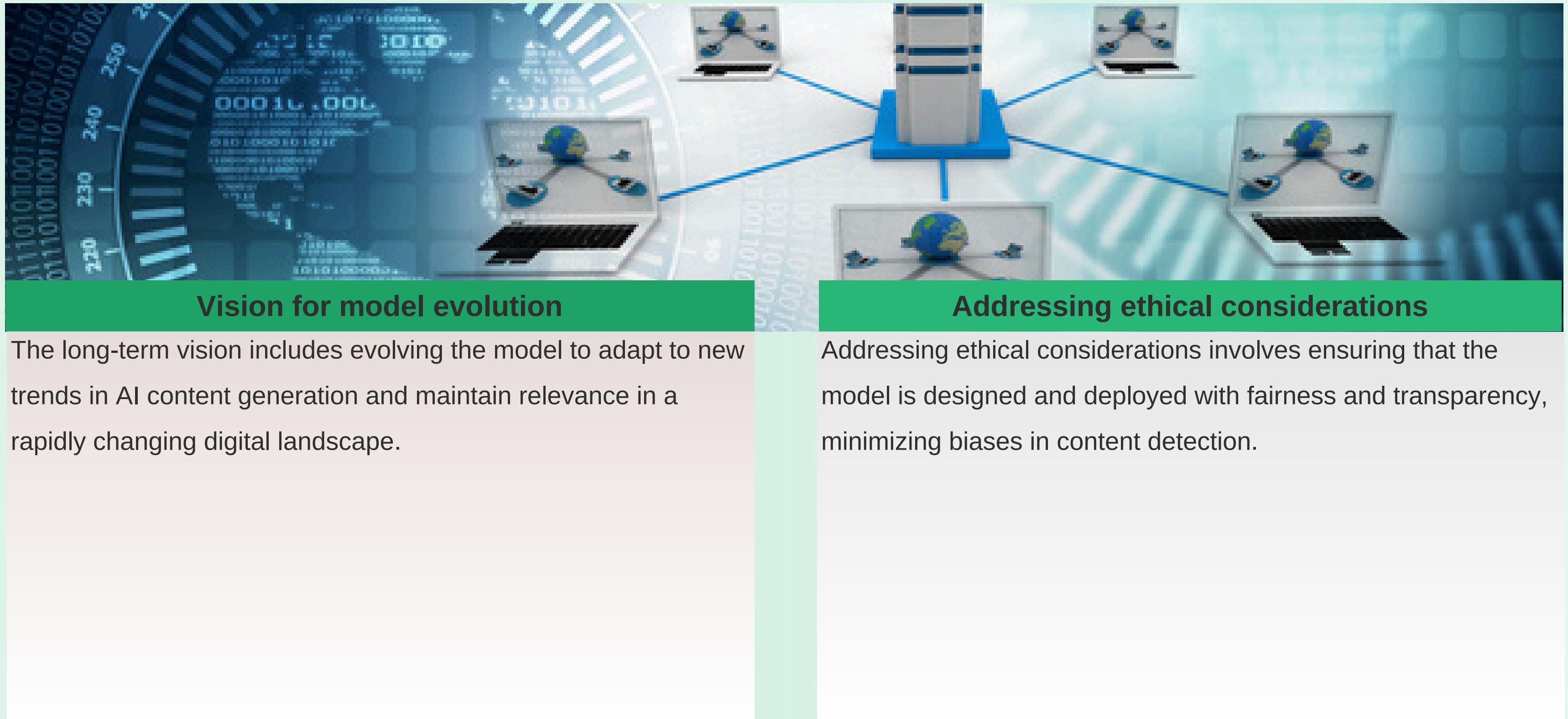
Mitigation strategies involve employing techniques like data augmentation and regularization methods to enhance robustness and prevent model overfitting.

# 06

## Future Directions



# Long-term Goals



## Conclusion:

The AI Content Detector represents a crucial advancement in verifying the authenticity of information. By effectively distinguishing between human and AI-generated content, we can uphold academic integrity, enhance journalistic credibility, and improve online discourse.

As we navigate an increasingly digital world, tools like this are essential for combating misinformation and fostering trust. Thank you for your attention, and we look forward to shaping the future of content verification together!



Thank you