# BiasCheck: A Political Bias Detector using RoBERTa

A PROJECT REPORT

Submitted by

Pranav Kaul [RA2211028010078]

Pranay Sharma [RA2211028010067]


*Under the Guidance of*

Dr. Saravanan M

Assistant Professor, Networking and Communications Department

*in partial fulfillments of the requirements for the degree of*


BACHELOR OF TECHNOLOGY

in

COMPUTER SCIENCE ENGINEERING

DEPARTMENT OF NETWORKING AND COMMUNICATIONS

COLLEGE OF ENGINEERING AND TECHNOLOGY

SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

KATTANKULATHUR- 603 203

MAY 2025

Department of Networking and Communications

**SRM Institute of Science & Technology**

**Own Work\* Declaration Form**

This sheet must be filled in (each box ticked to show that the condition has been met). It must be signed and dated along with your student registration number and included with all assignments you submit – work will not be marked unless this is done.

<u>To be completed by the student for all assessments</u>

**Degree/ Course**              :  **Btech Computer science and engineering with specialisation in Cloud Computing**

**Student Name**              : **Pranav Kaul, Pranay Sharma**

**Registration Number**      : **RA2211028010078, RA2211028010067**

**Title of Work**              : **Real-Time Political Bias Detection and Highlighting System for Web Content**

I / We hereby certify that this assessment compiles with the University's Rules and Regulations relating to Academic misconduct and plagiarism\*\*, as listed in the University Website, Regulations, and the Education Committee guidelines.

I / We confirm that all the work contained in this assessment is my / our own except where indicated, and that I / We have met the following conditions:

- Clearly referenced / listed all sources as appropriate

- Referenced and put in inverted commas all quoted text (from books, web, etc)

- Given the sources of all pictures, data etc. that are not my own

- Not made any use of the report(s) or essay(s) of any other student(s) either past or present

- Acknowledged in appropriate places any help that I have received from others (e.g. fellow students, technicians, statisticians, external sources)
- Compiled with any other plagiarism criteria specified in the Course handbook / University website

I understand that any false claim for this work will be penalized in accordance with the University policies and regulations.

| DECLARATION: |
| --- |
| I am aware of and understand the University's policy on Academic misconduct and plagiarism and I certify that this assessment is my / our own work, except where indicated by referring, and that I have followed the good academic practices noted above. |
| |
| If you are working in a group, please write your registration numbers and sign with the date for every student in your group. |

SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

KATTANKULATHUR – 603 203

BONAFIDE CERTIFICATE

Certified that 21CSP302L - Project report titled "**BiasCheck: A Political Bias Detector using RoBERTa**" is the bonafide work of "**Pranav Kaul RA2211028010078, Pranay Shrama RA2211028010067"** who carried out the project work under my supervision. Certified further, that to the best of my knowledge the work reported herein does not form any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

**SIGNATURE**

**Dr. Sarvanan M**

**SUPERVISOR**
**PROFESSOR**
DEPARTMENT OF
NETWORKING   AND
COMMUNICATIONS

**SIGNATURE**

**Dr. M Lakshmi**

**PROFESSOR &HEAD**
DEPARTMENT OF
NETWORKING       AND
COMMUNICATIONS

Examiner 1

Examiner 2

# ACKNOWLEDGEMENTS

# ABSTRACT

BiasCheck is an advanced political bias detection tool built with cutting-edge natural language processing (NLP) techniques, specifically leveraging the powerful RoBERTa (Robustly Optimized BERT Pretraining Approach) model. RoBERTa is known for its exceptional ability to understand and interpret the context in textual data. In today's fast-paced digital world, where misinformation and media biases spread quickly, BiasCheck provides a much-needed solution by helping users identify and understand political leanings in online content with speed and accuracy.

At its core, BiasCheck uses a variety of extensive datasets from trusted sources like AllSides, curated Kaggle datasets, and handpicked textual data, covering a wide range of political viewpoints. This allows the model to classify content into three categories: Left, Center, or Right, helping users spot both subtle and obvious biases in the media. By doing so, BiasCheck empowers users to critically evaluate digital media with more awareness and thoughtfulness.The development process of BiasCheck involved fine-tuning the RoBERTa model to enhance its ability to detect linguistic clues that signal political bias. After thorough training, the model underwent rigorous testing using accuracy and F1-score metrics, ensuring its reliability, precision, and consistency. The results from these tests have proven that BiasCheck performs robustly, providing accurate political bias detection across various types of text. What sets BiasCheck apart is not only its technical strength but also its accessibility. Delivered as a user-friendly Chrome extension, it integrates smoothly into daily browsing habits, offering real-time bias detection without interrupting the online experience. With an easy installation and straightforward operation, anyone can use BiasCheck, regardless of their technical skills. The broader impact of BiasCheck is significant, going beyond individual use. By helping users spot biased content instantly, it promotes a culture of media literacy and critical thinking. It encourages more informed conversations, supports transparency in media consumption, and plays an important role in reducing polarization and misinformation in society.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# ABBREVIATION

| ABBREVIATION | FULL FORM |
|---|---|
| NLP | Natural Language Processing |
| RoBERTa | Robustly Optimised BERT pretraining Approach |
| API | Application Programming Interface |
| DOM | Document Object Model |
| UI | User Interface |
| JSON | JavaScript Object Notation |
| GPU | Graphics Processing Unit |
| CORS | Cross-Origin Resource Sharing |
| JWT | JSON Web Token |
| LIME | Local Interpretable Model-agnostic Explanations |

# CHAPTER 1

# INTRODUCTION

Political bias in digital news distorts public perception, a problem exacerbated by online media. Automated NLP/ML solutions, like transformer models, can analyze news to detect bias, empowering readers and promoting transparency. This project develops a real-time detection system with a browser extension to provide immediate user feedback, addressing both technology and societal needs.

## 1.1 Introduction to BiasCheck :

In today's world, the rapid spread of information—both accurate and misleading—has made detecting bias in digital content more crucial than ever. Political bias, in particular, has the power to subtly sway public opinion, influence narratives, and deepen societal divisions. Traditional methods for identifying such bias, such as manual editorial reviews or crowd-based ratings, are often slow, inconsistent, and vulnerable to their own biases.

BiasCheck was created to tackle this issue, utilizing advancements in natural language processing to automate and standardize the detection of political bias in text. At its core, BiasCheck leverages RoBERTa, a transformer-based language model known for its contextual understanding and flexibility. The system is trained on a wide range of political articles, opinion pieces, and social media posts, all meticulously labeled to represent a diverse range of political perspectives. This allows BiasCheck to go beyond simple keyword detection, capturing complex patterns and linguistic indicators that point to underlying political leanings. The project follows a comprehensive, end-to-end workflow: starting with data collection and preprocessing, progressing through model training and validation, and culminating in real-time deployment via a Chrome extension. This extension empowers users to instantly analyze web content, providing clear visual cues and confidence scores regarding detected bias. Designed with accessibility in mind, BiasCheck ensures that individuals with different technical expertise can easily benefit from the tool. BiasCheck's impact goes beyond individual users. By offering the public reliable tools for detecting bias, the project

aims to raise the bar for media literacy, promote critical engagement with information, and encourage healthier, more transparent public discourse. In doing so, it addresses a critical gap in the digital information ecosystem and sets a new standard for responsible, AI-driven media analysis.

## 1.2 Motivation

The motivation behind BiasCheck stems from the growing awareness that subjective and political biases are deeply woven into the language of digital media. As society becomes more diverse and interconnected, the consequences of unchecked bias—such as marginalization, misinformation, and heightened polarization—are becoming more pronounced. While there is widespread recognition of these dangers, effective and accessible tools to assess the objectivity of daily content remain scarce.

Existing NLP solutions have primarily focused on detecting toxic language or explicit hate speech, leaving the more complex challenge of identifying subjective and political bias largely unexplored. BiasCheck seeks to fill this gap by offering a practical, automated tool that helps writers, readers, and communication professionals spot and mitigate bias in real-time. The development of the tool is motivated by a desire to encourage more objective, inclusive, and transparent communication in both professional and public spheres.

The technical backbone of BiasCheck—fine-tuning RoBERTa on carefully selected datasets—was chosen for its ability to understand contextual nuances and subtle linguistic cues. The project also acknowledges the need for continual improvement, incorporating user feedback and analytics to enhance its performance over time. By offering bias detection through a simple browser extension, BiasCheck makes critical media engagement more accessible, enabling users to be more thoughtful consumers and creators of information.
At its core, the motivation for BiasCheck is not solely technological but also societal: to promote fairness, reduce polarization, and foster a more informed and resilient public that can navigate the complexities of today's information landscape.

**1.3 Sustainable Development Goal of the Project**

Sustainable Development Goal 16 (SDG 16), which focuses on promoting peaceful and inclusive societies, ensuring access to justice for all, and building effective, accountable, and inclusive institutions, is highly relevant to the BiasCheck project. In today's digital landscape, the rapid spread of information—often colored by political bias—can undermine trust in media and institutions, fuel polarization, and erode democratic processes. BiasCheck directly addresses these challenges by equipping users with a tool that detects and highlights political bias in digital content, thereby fostering greater transparency and accountability in the information ecosystem. By making it easier for individuals to identify and critically assess biased narratives, BiasCheck supports the core SDG 16 objective of ensuring public access to reliable and impartial information. This not only empowers citizens to make informed decisions but also encourages media organizations to uphold higher standards of objectivity and fairness, contributing to the development of more trustworthy and inclusive institutions. Furthermore, by promoting media literacy and critical thinking, BiasCheck helps mitigate the risks of misinformation and societal division, aligning closely with SDG 16's vision of peaceful, just, and resilient societies built on the foundation of transparent and accountable information systems.

**1.4 Product Vision Statement**

**1.4.1 Audience:**

○ **Primary Audience:** Everyday news readers, students, educators, and journalists seeking to understand and critically evaluate the political bias present in digital content.

○ **Secondary Audience:** Media organizations, fact-checkers, researchers, and policy analysts interested in monitoring, benchmarking, or improving the objectivity of public communications

### 1.4.2 Needs:

○ **Primary Needs:** Real-time, accurate detection of political bias (Left, Center, Right) in news articles, opinion pieces, and social media posts.

○ Intuitive browser integration for seamless, non-intrusive analysis during everyday reading.

○ Clear visual cues and confidence scores to support critical media consumption and decision-making.

○ **Secondary Needs:**

○ Summarized bias reports and exportable analytics for educational, editorial, or research purposes.

○ Feedback channels for users to contribute to ongoing model improvement.

○ Compatibility with diverse content formats and evolving digital platforms.

### 1.4.3 Products:

○ **Core Product:** A Chrome extension that leverages a fine-tuned RoBERTa model to analyze and highlight political bias in digital text, providing instant feedback within the user's browsing experience.

○ **Additional Features:**

○ Visual bias highlighting (color-coded as Left, Center, or Right) and summary panels with confidence scores.

○ Exportable reports (CSV/PDF) for classroom, newsroom, or research use.

○ User feedback integration for continuous model refinement and adaptation to new discourse trends.

### 1.4.4 Values:

- **Core Values:**

- **Objectivity:** Rigorous model training on diverse, reputable datasets to ensure fair and balanced bias detection.

- **Transparency:** Clear explanations of bias ratings and accessible audit trails for users and stakeholders.

- **Empowerment:** Promoting media literacy and critical thinking by making bias detection accessible to all.

- **Differentiators:**

- **Advanced NLP:** Utilizes state-of-the-art RoBERTa architecture for nuanced, context-aware bias detection.

- **User-Centric Design:** Simple, intuitive interface requiring no technical expertise.

- **Continuous Improvement:** Incorporates user feedback and analytics for ongoing accuracy and relevance in a dynamic media landscape.

## 1.5 Product Goal

The primary goal of BiasCheck is to help individuals become more discerning and informed consumers of digital content by offering an intuitive, reliable, and accessible tool for detecting political bias in real-time. In an age where misinformation and polarized narratives can easily sway public opinion and undermine democratic values, BiasCheck aims to bridge the gap between technology and media literacy. By harnessing the advanced contextual understanding of the RoBERTa language model, the tool seeks to accurately identify and classify political leanings—Left, Center, or Right—in a wide variety of textual content, from news articles to social media posts. This allows users to quickly evaluate the objectivity of the information they encounter, helping them make more balanced decisions about what to trust and share.

BiasCheck's user-friendly Chrome extension is designed to integrate seamlessly into everyday browsing, ensuring that bias detection is both effortless and unobtrusive. The tool not only highlights biased content but also provides confidence scores and source metadata, offering users deeper insights into the perspectives underlying the information they consume. Beyond individual empowerment, BiasCheck aims to raise overall media standards by promoting greater transparency and accountability among content creators and publishers. By encouraging critical engagement and responsible information consumption, BiasCheck contributes to a healthier, more informed, and less polarized digital society. In essence, BiasCheck strives to make bias detection an accessible part of daily digital life, supporting both personal empowerment and the broader goal of building trustworthy, transparent information ecosystems.

**1.6 Product Backlog**

| S.No | User Stories of BiasCheck |
|---|---|
| #US 1 | As a **casual reader**, I want to run a quick bias scan on any news article so that I can gauge its political slant. |
| #US 2 | As a **student**, I want to see highlighted passages classified by bias so that I can learn how subtle language shapes opinion. |
| #US 3 | As a **journalist**, I want to view bias statistics for my own drafts so that I can ensure my reporting stays balanced. |
| #US 4 | As a **researcher**, I want to export bulk bias-analysis results in CSV/PDF so that I can perform in-depth academic studies.. |
| #US 5 | As a **developer**, I want bias-thresholds and label-mapping to be configurable via a settings file so that I can adapt BiasCheck to different domains. |
| #US 6 | As an **API consumer**, I want a REST endpoint that returns bias labels and confidence scores so that I can integrate BiasCheck into other tools. |
| #US 7 | As a **policymaker**, I want to track bias trends over time across multiple outlets so that I can monitor shifts in media framing. |
| #US 8 | As an **educator**, I want to demonstrate bias-detection live in my classroom via the Chrome extension so that students can interactively learn media literacy. |
| #US 9 | As a **browser user**, I want real-time bias notifications when I open an article so that I'm immediately aware of any leaning. |

| #US 10 | As a **compliance officer**, I want to log all bias-analysis events with timestamps so that I can audit content review processes. |
|---|---|
| #US 11 | As an **end-user**, I want access to a built-in glossary of bias terminology so that I can understand each label and its implications. |

Table 1.1 Shows different user stories for biascheck usage

The product backlog of the BiasCheck platform was configured using the MS Planner Agile Board, as illustrated in Figure 1.1 below. This backlog comprises the comprehensive set of user stories for BiasCheck, each articulated with MoSCoW prioritization, clearly defined functional and non-functional requirements, and detailed acceptance criteria linked to corresponding development tasks.



Figure 1.1 MS Planner Board of BiasCheck Extension

## 1.7 Product Release Plan

The following Figure 1.2 depicts the release plan of the project



**Release Plan/ Roadmap**

| | W/1 | W/2 | W/3 | W/4 | W/1 | W/2 | W/3 | w/4 |
|---|---|---|---|---|---|---|---|---|
| Proposal | Define Product Vision and Finalise Datasets | | | | | | | |
| Feature # 1 | | Fine Tune BERT using labelled biased datasets | | | | | | |
| Feature # 2 | | Create FrontEnd for the chrome extension | | | | | | |
| Feature # 3 | | | | Create API using Flask and establish endpoints | | | | |
| Feature # 4 | | | | | | Test bias classification logic and connect frontend and backend | | |
| Feature # 5 | | | | | | Refine Model, Fix Bugs, Collect Feedback | | |

Figure 1.2 Release Planner sheet of BiasCheck Extension

# CHAPTER 2

# SPRINT PLANNING AND EXECUTION

The report documents the systematic approach undertaken to develop the political bias detection system. It details the division of work into focused sprints, describing the key activities, challenges, and milestones achieved in each phase. This chapter highlights the collaborative efforts, technical decisions, and iterative improvements made during dataset collection, model fine-tuning, backend deployment, and frontend integration. It also reflects on lessons learned and adjustments implemented to ensure timely delivery and quality outcomes.

## 2.1 Sprint 1

### 2.1.1 Sprint Goal with User Stories of Sprint 1

In the first sprint, the team set out to establish the foundational components of BiasCheck, focusing on creating a minimum viable product that could ingest text, perform basic bias classification with RoBERTa, and display results through a simple Chrome extension. Key user stories included setting up the data collection and preprocessing pipeline, fine-tuning the RoBERTa model on a small, representative dataset, exposing an API endpoint for bias inference, and implementing a rudimentary UI for highlighting bias directly on web pages. By the end of the two-week period, these core capabilities formed the backbone of our bias detection system.

| S.NO | Detailed User Stories |
|------|----------------------|
| US #1 | As a student, I want to see highlighted passages classified by bias so that I can learn how subtle language shapes opinion. |
| US #2 | As a journalist, I want to view bias statistics for my own drafts so that I can ensure my reporting stays balanced. |
| US #3 | As a researcher, I want to export bulk bias-analysis results in CSV/PDF so that I can perform in-depth academic studies. |

Table 2.1 Detailed User Stories of sprint 1

Figure 2.1 user story for building the extension



Figure 2.2 user story for understanding the political bias and awareness



Figure 2.3 User story for highlighting text

**2.1.2 Functional Document**

**2.1.2.1. Introduction**

The Functional Document for Sprint 1 provides a detailed exposition of the system's capabilities, interfaces, and constraints as implemented in this initial development cycle. It delineates how user requirements translate into concrete services, how data flows through the system, and how the user interacts with the application. By codifying inputs, outputs, and processing logic, this document ensures that all team members share a consistent understanding of the core functionality delivered in Sprint 1

**2.1.2.2. Product Goal**

The primary goal of Sprint 1's functional specification is to enable BiasCheck to perform basic political bias detection end-to-end. This involves accepting raw text input, normalizing it for compatibility with RoBERTa, executing inference to assign a bias label and confidence score, and returning results in a standardized format. The functional scope also covers inline highlighting of biased text within the browser extension. These capabilities form the minimum viable product, demonstrating the feasibility of real-time bias detection and guiding subsequent iterations.

**2.1.2.3. Demography (Users and Deployment Regions)**

BiasCheck's Sprint 1 features are intended for a diverse user base, including casual readers seeking quick bias checks, students analyzing source credibility, and journalists verifying editorial neutrality. These users operate within standard desktop environments, using Chrome as their primary web browser. On the development side, engineers and QA analysts access the backend API locally or via a staging server. The system assumes a stable internet connection for model hosting, and that users possess basic familiarity with browser extensions and modern web interfaces.

**2.1.2.4 Core Operational Modules Key operational flows include:**

The functional workflow begins when a user selects text on any web page and clicks the BiasCheck icon. The extension captures the selected text, wraps it in a JSON payload, and sends it over HTTPS to the backend inference API. Upon receipt, the

API preprocesses the text—performing tokenization, lowercasing, and removal of special characters—before feeding it into the fine-tuned RoBERTa model. The model returns a bias label (Left, Center, or Right) and a numeric confidence score. The API then packages these results into a JSON response. Back in the extension, the text is annotated inline with color-coded highlights corresponding to the detected bias, and a summary panel displays the overall bias classification and score.

**2.1.2.5 Key Platform Features:**

In Sprint 1, the application delivered a suite of core features that collectively enable basic bias detection functionality. First, the text ingestion and preprocessing component standardizes incoming content by normalizing whitespace, stripping special characters, and tokenizing sentences to prepare inputs for the RoBERTa model. Once text is preprocessed, the RoBERTa inference API accepts these inputs through a RESTful Flask endpoint; it processes the data and returns a bias classification along with a confidence score. On the frontend side, the inline highlighting extension seamlessly injects color-coded markers into the Document Object Model (DOM), visually indicating detected bias as users browse articles. A collapsible result panel complements this visual feedback by summarizing the overall bias label, displaying the confidence percentage, and providing metadata about the data source. Finally, built-in error handling safeguards the user experience by validating requests and offering clear feedback whenever inputs are empty or malformed, preventing silent failures and guiding corrective action.

**2.1.2.6 Authorization Matrix**

| Role | Access Level |
|---|---|
| Developer | Full access to codebase, model configurations, and API |
| QA Analyst | Permission to execute and validate API endpoints and UI |
| End User | Ability to invoke text analysis and view results via UI |

Table 2.2 Representing the authorization matrix

**2.1.2.7 Assumptions**

Several assumptions underpin Sprint 1's functionality: the RoBERTa model is pre-trained and accessible within the container environment; users select sufficient text segments (at least 50 tokens) to yield meaningful bias predictions; the hosting environment supports Python 3.8+, Flask, and required NLP libraries; and the Chrome browser is version 80 or higher to ensure compatibility with extension APIs. Future sprints will revisit and refine these assumptions as the system scales.

Development began with establishing a Git repository and CI pipeline. The data team authored preprocessing scripts to normalize text, remove noise, and convert content into RoBERTa's input format. Simultaneously, the ML team fine-tuned a pre-trained RoBERTa checkpoint on a small, representative sample from AllSides and Kaggle sources. The containerized model was served via a Flask-based API, and the frontend team scaffolded the extension to listen for user selection events and apply bias highlights dynamically.

**2.1.3 Architecture Document**

**2.1.3.1. Application**

The BiasCheck application employs a hybrid architecture with a client-side Chrome extension and a server-side inference service. The extension features three JavaScript modules: a Content Script that extracts text, highlights biased sections using CSS overlay via shadow DOM, and responds to user interactions; a Background Script managing the extension's lifecycle, user preferences in Chrome's Local Storage, network checks, and inter-module communication; and a React-based Popup UI (bundled with Webpack) for configuration, full-page scans, and displaying summarized analysis results.

The server-side inference service, containerized with Docker and running a Python Flask application, exposes an /api/v1/analyze REST endpoint. It receives text, preprocesses it using a Preprocessing Module (whitespace, unicode, RoBERTa tokenizer), and feeds it to a fine-tuned RoBERTa model (leveraging GPU if available). Predictions (bias labels and confidence scores) are returned as JSON with model metadata. Logging middleware tracks inference calls. Configuration (model version, input limits, concurrency) is managed via environment variables and a mounted configuration file, while Docker secrets handle sensitive API keys. This client-server separation and containerization support independent scaling and resilience.

This application architecture describes how the Chrome extension interfaces with the bias detection engine and supporting modules. Table 2.3 lists the major components—Content Script, Background Script, Preprocessing Module, and Model Inference Engine—along with their responsibilities. As depicted in Figure 2.1, the Content Script extracts selected text from web pages and communicates with the Inference Engine either directly in-browser (via ONNX.js or TensorFlow.js) or through an optional local Flask API. The Background Script manages lifecycle events and user preferences stored in Chrome Local Storage, while the Popup Interface provides controls for extension settings and displays analysis summaries.

| Component | Description |
|-----------|-------------|
| Chrome Extension Frontend | UI component embedded into web pages, includes popup and content script. |
| Background Script | Manages lifecycle events and background processing of the extension. |
| Content Script | Injected into web pages; extracts text and highlights biased content. |
| Model Inference Engine | Loads the fine-tuned BERT/Roberta model (via TensorFlow.js or ONNX.js) in-browser or sends to a local Flask/Node API. |
| Local API (optional) | Flask or Node.js server to host the model for inference if the browser is too slow. |
| Preprocessing Module | Cleans and tokenized text before inference. |

Table 2.3 Different components of the extension

## 2.1.3.2 System Architecture-

The overall system architecture for BiasCheck adheres to modern best practices for scalable, resilient web applications. At its core lies a decoupled, three-tier design: the Presentation Layer (Chrome extension), the Application Layer (Flask inference service), and the Infrastructure Layer (container orchestration and observability). The Presentation Layer runs entirely within the user's browser, minimizing latency for user interactions and reducing server load by offloading simple UI tasks. The Application Layer is provisioned via container orchestration tools (Docker Compose in development, Kubernetes in production), allowing horizontal scaling of inference replicas behind a load balancer.

Within the Application Layer, the service is configured to run multiple worker processes (using Gunicorn) to handle concurrent inference requests. Each worker maintains a persistent model instance in memory for rapid inference, reducing cold-start overhead. A reverse proxy (NGINX) sits in front of the Flask application to manage SSL termination, implement gzip compression, and enforce HTTP/2 for

efficient multiplexing. Rate limiting and IP-based access controls are enforced at the proxy level to mitigate abusive or malicious traffic. For environments requiring high availability, the service can be deployed across multiple availability zones with health checks and automatic failover configured via the orchestration platform.

Observability is embedded throughout the system: detailed request and response logs are pushed to a centralized logging service (e.g., Elasticsearch) via Filebeat agents, while metrics—such as request rates, latency histograms, and error percentages—are emitted to a monitoring system (e.g., Prometheus) with Grafana dashboards for real-time visualization. Distributed tracing is supported through OpenTelemetry integrations in both the Flask app and the extension's network layer, enabling end-to-end latency analysis from the browser event to the model inference.

Security and compliance are top priorities. All data in transit is encrypted with TLS 1.3, and Cross-Origin Resource Sharing (CORS) policies are narrowly scoped to the extension's origin to prevent unauthorized cross-site requests. Authentication tokens (JWTs) can be integrated for future versions requiring user-specific analytics or protected model versions. Infrastructure components are hardened according to CIS benchmarks, with automated vulnerability scanning of container images incorporated into the CI/CD pipeline. Deployment manifests define resource limits and readiness probes to ensure stable rollouts and rapid rollback in case of anomalies.



Figure 2.4 System Architecture Diagram

**2.1.3.3. Data Exchange Contract:**

BiasCheck operates as a hybrid system, integrating a client-side Chrome extension with a server-side inference service to detect bias within a user's browsing experience. The Chrome extension consists of a Content Script, a Background Script, and a Popup UI, each handling specific tasks to ensure seamless bias detection. The Content Script extracts and highlights potentially biased text on web pages. The Background Script manages the extension's lifecycle, user preferences, network connectivity, and inter-module communication. The React-based Popup UI allows users to configure settings and view analysis results.

The server-side inference service, a Python Flask application within a Docker container, exposes a RESTful API endpoint (/api/v1/analyze). It receives text, preprocesses it, and uses a fine-tuned RoBERTa model to predict bias labels and confidence scores, returning this data as a JSON response. Inference requests and responses are also asynchronously published to a message queue (like RabbitMQ or AWS SQS) for usage analytics, which are stored in a time-series database (e.g., InfluxDB) to drive dashboards and model retraining efforts.

For large data operations, such as importing training data, BiasCheck utilizes file-based exchange via signed URLs to cloud storage (e.g., AWS S3), with a separate "ingest" API call to initiate processing, thus avoiding overloading the inference service.

Client-side data management involves Chrome Local Storage for immediate access to user preferences and IndexedDB for caching recent analysis results per URL, significantly reducing redundant API calls.

Critical system data is centralized in PostgreSQL (model metadata, user consents, retraining jobs), while unstructured data (raw text, logs, annotations) resides in MongoDB, facilitating flexible analytics. All inter-service communication is secured with TLS 1.3, and access is controlled through scoped service accounts and role-based permissions.

While future enhancements may include WebSocket-based live collaboration, the initial architecture leverages REST, message queues, and client-side storage to establish a secure, scalable, and responsive Data Exchange Contract.

| Source | Destination | Data | Frequency | Mode |
|--------|-------------|------|-----------|------|
| Content Script | Inference Module | Text snippets from web pages | On page load / scroll | Function call / API |
| Inference Module | Content Script | Bias label, confidence score | Real-time | JS Function return / Fetch |
| Popup Script | Chrome Storage | User preferences (e.g., bias threshold) | On change | Chrome API |

Table 2.4 Representing data exchange

## 2.1.4 UI DESIGN



Figure 2.5.1 BiasCheck Highlighting different biases

Figure 2.5.2 Chrome Extension in Developer mode



Figure 2.6 BiasCheck available in Chrome extensions

19

## 2.1.5 Functional Test Cases

| Feature | Test Case | Steps to Execute Test Case | Expected Output | Actual Output | Status | More Information |
|---------|-----------|----------------------------|-----------------|---------------|--------|-----------------|
| Bias Detection API | Valid Text Classification | 1. Input a politically biased sentence into the API. 2. Submit the request. | The API should return a valid political bias classification label (Left, Center, Right) with a confidence score. | API returns classification as "Left" with 87% confidence. | Pass | Confirm that the classification is logical. Test multiple sample inputs across different biases to verify accuracy. |
| Chrome Extension | Highlight Detected Bias | 1. Open a webpage with political content. 2. Activate the extension. | Detected biased sentences should be highlighted with color codes: - Red for Right - Blue for Left - Grey for Center | Text is highlighted correctly with corresponding colors. | Pass | Hovering over text shows tooltip with bias and confidence score. |
| Chrome Extension | DOM Manipulation | 1. Run extension | Page layout remains | Layout remains | Pass | Ensure it doesn't |

| | | | | | |
|---|---|---|---|---|---|
| | | on a website with dynamic content (e.g., news site with infinite scroll). 2. Observe page structure before and after extension activates. | intact. Only text color and tooltip are added. | unchanged, and only text is altered with highlight. | | break site CSS or interfere with page functionality. |
| Real-Time Detection | Content Change Monitoring | 1. Scroll a news feed (e.g., Twitter). 2. New content loads. 3. Bias detection should run on new content. | Newly loaded content is automatically detected and highlighted without refreshing or manual intervention. | New content is detected, and bias is highlighted dynamically. | Pass | Useful for infinite scrolling content platforms. |
| Browser Compatibilit y | Chrome Compatibility Check | 1. Install the extension on latest | Extension activates without errors and | Works perfectly on Chrome Version 123. | Pass | Future test needed for Edge and Chromium-b |

| | | Chrome version. 2. Open multiple news sites and activate. | bias detection works as expected. | | | ased browsers. |
|---|---|---|---|---|---|---|
| Extension UI | Toggle Activation Button | 1. Click extension icon to turn bias detection on/off. 2. Observe changes on page. | Bias detection should start when toggled ON and all highlights should be removed when toggled OFF. | Toggle works, text is highlighted on ON and cleaned on OFF. | Pass | Add animation or confirmation for toggle state (optional). |
| User Feedback (Future) | Submit Feedback on Misclassified Text | 1. Click on a highlighted sentence. 2. Select "Report Incorrect Bias". 3. Enter preferred label and submit. | System should store the feedback locally or send to backend for model training. | Not implemented (Planned for future sprint). | N/A | Will enhance long-term model accuracy if feedback loop is integrated. |

Table 2.5 Detailed Functional Test Case

## 2.1.6 Sprint Retrospective



| Sprint Retrospective 1 | | | |
|---|---|---|---|
| What went well | What went poorly | What ideas do you have | How should we take action |
| Bias detection algorithm returned consistent and accurate results across sample inputs. | Extension had issues detecting bias in dynamically loaded content (e.g., infinite scroll). | Integrate MutationObserver to monitor dynamic DOM changes. | Implement and test MutationObserver in the next sprint. |
| Chrome extension injected styles and tooltips correctly without breaking the page layout. | Confidence scores were inconsistent for longer paragraphs. | Implement paragraph-level averaging for long-form content. | Adjust API threshold or post-process confidence scores for consistency. |
| CI pipeline executed automated tests successfully | Test coverage lacked edge case scenarios for preprocessing | Expand unit tests to include edge cases in text preprocessing | Add additional test cases and increase coverage metrics |

Figure 2.7  Sprint Retrospective for the Sprint 1

## 2.2 SPRINT 2

### 2.2.1 Introduction

Sprint 2's Functional Document details how BiasCheck's new customization and in-context learning features are integrated into the existing bias-analysis pipeline. This document outlines the Settings Module, which provides APIs enabling users to adjust the bias-sensitivity threshold and toggle glossary tooltips. It also describes the Glossary Service, responsible for delivering term definitions upon request, and the enhancements made to the inference flow to incorporate these user preferences in real-time. This specification aims to align all frontend, background, and backend components, ensuring a clear and user-driven workflow for the BiasCheck Chrome extension, which highlights content on webpages based on its identified political bias (red, yellow, and blue for left, neutral, and right-wing biases, respectively).

### 2.2.2 Functional Document

### 2.2.1.1 Product Goal

By the end of this sprint, BiasCheck will offer users comprehensive control over its detection behavior and provide immediate educational context. Specifically, users will be able to fine-tune the "political slant" sensitivity using a slider, enable or disable glossary pop-ups for any highlighted term, and rely on their chosen settings to persist across browsing sessions.

These enhancements build upon the core bias-scanning engine developed in Sprint 1, transforming BiasCheck into a more personalized and learning-oriented browser extension.

**2.2.1.2 Demography (Users and Context)**

The intended users of BiasCheck can be grouped into several categories. Casual readers are looking for quick, easily accessible cues about potential bias and simple definitions for unfamiliar terms. They typically use Chrome on desktop or mobile devices and expect immediate feedback from the extension as they browse web pages. Researchers and students, on the other hand, require more precise control over bias detection, with the need for accurate threshold tuning and reliable glossary lookups for deeper analysis. These users often have multiple tabs open and expect consistent performance from BiasCheck across browser restarts. Lastly, developers and QA personnel test the extension in staging and local development environments. They use Chrome v80 or later, along with Node/Flask services to simulate the Settings and Glossary APIs, ensuring that new features work properly before release.

**2.2.1.3 Business Processes**

Several key business processes are involved in delivering the new features. Settings Management occurs when a user adjusts the bias slider or the glossary toggle within the extension's popup UI. The extension's background script then interacts with a Settings API, which records the new values in Chrome Local Storage. Subsequently, the background script broadcasts an update to all active content scripts, ensuring that the changes take effect immediately on all open pages. Glossary Lookup is initiated when a user hovers their cursor over a highlighted term. The content script responds to this action by fetching the term's definition from the Glossary Service. In cases where the Glossary Service is temporarily unavailable, the content script uses a local cache to provide a fallback definition or displays an "Unavailable" message to the user. Bias Inference with Preferences describes how the extension uses the user's settings during bias analysis. Before sending a text segment to the RoBERTa engine for classification, the content script retrieves the user's current bias threshold. This threshold is included in the payload sent to the engine, ensuring that the bias classification respects the user's desired sensitivity level. Persistence & Sync refers to how user preferences are maintained across browsing sessions. When the extension loads, or after Chrome restarts, the background script reads the saved preferences from Chrome Local

Storage. It then communicates these preferences to all open tabs, so the user's configuration is applied automatically without requiring manual page reloads.

## 2.2.1.4 Features

Sprint 2 introduces three primary features to the BiasCheck extension. First, the Settings API provides methods to retrieve and modify bias thresholds, represented as a floating-point value between 0.0 and 1.0, and to enable or disable glossary tooltips. These settings are applied immediately across all web pages the user visits. Second, the Glossary Service is implemented as a lightweight REST endpoint. This service delivers term-definition pairs from a versioned JSON data store, with each definition lookup designed to complete within 200 milliseconds. Third, Preference Persistence ensures that any user-defined settings are stored in Chrome Local Storage. These settings are automatically reapplied whenever the browser or the extension is restarted, providing a seamless and consistent user experience.

## 2.2.1.5 Authorization Matrix

| Role | Access & Capabilities |
|---|---|
| Developer | Full access to source code, Settings API schema, Glossary JSON, and integration tests |
| QA Analyst | Can execute end-to-end scenarios: adjust settings, trigger glossary lookups, & verify persistence |
| End User | May adjust thresholds via UI, hover for definitions, and view bias highlights |

Table 2.6 Roles and Access

### 2.2.1.6 Assumptions

Several assumptions underpin the design and implementation of these features. It is assumed that Chrome's storage quotas are sufficient for storing user preference data, which is expected to be less than 10 KB per user. The Glossary JSON file is assumed to remain smaller than 500 KB, which is necessary to meet the 200-millisecond lookup performance target. The extension is designed for users running Chrome version 80 or later, which supports the required ES6 JavaScript features. It is also assumed that the inference endpoint used by the extension accepts a runtime threshold parameter, allowing the extension to dynamically adjust bias detection sensitivity. Finally, it is assumed that users generally have network access, with local caching mechanisms in place to handle brief network outages.

### 2.2.3 Architecture Document

### 2.2.2.1 Application Architecture

In Sprint 2, BiasCheck's application architecture matured into a truly modular system in which the Chrome extension frontend, a high-performance bias-inference engine, and two newly introduced microservices—the Settings Service and the Glossary Service—operate in concert to deliver a seamless, personalized experience. The Settings Service lives alongside the RoBERTa inference endpoint in a Dockerized Flask container and exposes a REST API for persisting and retrieving user preferences. When the extension popup is opened—or immediately upon browser startup—the Background Script fetches the current settings (bias threshold and glossary-toggle flag) and broadcasts them to all active Content Scripts via Chrome's runtime messaging. Each Content Script then adapts its behavior in real time: inference requests include the up-to-date threshold, and hover listeners activate only if glossary lookups are enabled. Meanwhile, the Glossary Service draws from a version-controlled JSON store of political term definitions and answers lookup requests in under 200 ms. On hover, the Content Script asynchronously fetches the definition and renders it within a Shadow DOM tooltip, ensuring that page styles remain untouched. Comprehensive error handling and retry logic are built in: if either service is unreachable, cached data is used or a "Definition unavailable" notice appears, preserving functionality without degrading the user's browsing flow. Figure 2.3 (placeholder) illustrates the full component interaction,

showing how extension scripts, the Flask container, and the inference engine exchange messages to uphold both performance and modularity.

## 2.2.2.2 Storage Architecture

To balance persistence, performance, and offline resilience, BiasCheck uses a two-tiered browser-side storage strategy. Core user preferences are stored in Chrome Local Storage under a structured namespace (biascheck.preferences), consuming only a few kilobytes while persisting across sessions and synchronizing automatically when users are signed in to Chrome. Each preference set includes the numerical bias threshold and a boolean flag indicating whether glossary tooltips should appear. For richer educational content, the extension uses IndexedDB as a local cache for glossary entries. The first time a term is looked up, its definition and metadata are saved in an object store along with a timestamp. A background cleanup routine purges entries older than 24 hours, enforcing a time-to-live policy that ensures freshness while minimizing network calls. In situations of limited connectivity, the extension seamlessly falls back to these cached definitions, maintaining uninterrupted tooltips for recently viewed terms. Table 2.7 (placeholder) provides a detailed breakdown of each storage mechanism—its capacity, data schema, eviction policy, and primary use case—explaining how BiasCheck efficiently manages both user preferences and glossary data.

## 2.2.2.3 Data Exchange Contract

The Data Exchange Contract establishes a precise, versioned specification for all cross-component communications introduced in Sprint 2, defining payload schemas, triggering events, and stringent performance targets to guarantee a fluid user experience. Content Scripts send inference requests as JSON objects containing text (the selected article snippet), threshold (a float between 0.0 and 1.0), and an optional sessionId (UUID) to the RoBERTa endpoint; this endpoint must return { "label": "Left" | "Center" | "Right", "confidence": float } within a 500 ms SLA. Settings updates flow over Chrome's native messaging API in messages of the form { "type": "SETTINGS_UPDATE", "payload": { "threshold": float, "glossaryEnabled": bool } }, with a delivery budget of 50 ms to all open tabs. Glossary lookups issue { "term": string } to the Glossary Service and expect { "definition": string, "source": string } in under 200 ms. To maximize reliability, each channel implements exponential-backoff retry logic—up to two attempts on transient errors—before falling back to cached results or user-facing notifications. All interactions are instrumented

with telemetry hooks that log request timestamps, latencies, HTTP statuses, and schema-validation results. Table 2.8 (placeholder) enumerates each transaction—source, destination, payload structure, trigger condition, and maximum allowable latency—providing a definitive reference for both frontend and backend teams to ensure end-to-end consistency and responsiveness.

**2.2.4 UI Design**



Figure 2.8 Implementation UI

**2.2.5 Functional Test Cases**



| Functional Test Case Sprint 2 | | | | | | |
|---|---|---|---|---|---|---|
| Feature | Test Case | Steps to Execute Test Case | Expected Output | Actual Output | Status | More Information |
| Settings Panel | Adjust Bias Threshold Slider | 1. Open extension popup. 2. Move bias threshold slider to 0.75. | Analysis reflects the new threshold: content with bias score ≥0.75 is highlighted; lower | ighlights adjust as e> | Pass | Test with multiple threshold values to verify dynamic updates. |
| Inline Glossary Tooltips | Display Definition Tooltip | 1. Enable glossary in settings. 2. Highlight a political term on a webpage. 3. Hover over the term. | A tooltip appears within 200â€¯ms showing the correct definition and source link. | efinitions display cc | Pass | Ensure tooltip does not interfere with page styles. |
| Preference Persistence | Persist Settings Across Sessions | 1. Set threshold to 0.6 and disable glossary. 2. Restart the browser. | Threshold remains at 0.6 and glossary remains disabled without requiring manual | ttings persisted afte | Pass | Test on multiple restart cycles and across different profiles. |
| Responsive UI Layout | Adapt Panel for Mobile View | 1. Resize browser to 320Â—568 (mobile viewport). 2. Open extension popup. 3. Navigate through settings and close popup. | Settings panel stacks vertically, controls remain fully visible and touch targets ≥44â€¯px. | adapts without layou | Pass | erify using both browser dev tools and a physical mobile device |
| Offline Glossary Fallback | Handle Glossary Lookup When Offline | 1. Disable network connection. 2. Enable glossary in settings. 3. Hover over a previously cached term. | A fallback message 'Definition unavailable' appears, and no errors are thrown. | fline notice shown g | Pass | Test both cached and uncached terms for fallback behavior. |

Figure 2.9 showcasing the Function Test cases 2nd sprint

## 2.2.6 Sprint Retrospective

| Sprint Retrospective 2 | | | |
|---|---|---|---|
| **What went well** | **What went poorly** | **What ideas do you have** | **How should we take action** |
| User-configurable thresholds and glossary integration were implemented smoothly. | Frontend-backend JSON schema mismatches caused integration delays. | Standardize JSON contracts and enforce schema validation. | Document and enforce API schemas in the development guidelines. |
| UI design became more responsive and accessible. | Glossary caching strategy led to stale entries during initial tests. | Implement cache invalidation with expiration policies. | Add TTL (time-to-live) to cached glossary data and clear on updates. |
| Backlog grooming and estimation accuracy improved. | Performance degradation when loading large glossary files. | Lazy-load glossary entries to reduce initial load time. | Refactor glossary module to fetch entries on demand. |

Figure 2.10 showcasing the sprint retrospective of 2nd sprint

# CHAPTER 3

# RESULTS AND DISCUSSION

The report presents a comprehensive evaluation of the system's performance, including quantitative metrics and qualitative insights. It analyzes the effectiveness of the fine-tuned RoBERTa model in classifying political bias across diverse news articles and discusses the real-time highlighting mechanism's impact on user experience. This chapter also addresses observed limitations, challenges encountered during testing, and the implications of the findings for future research. The discussion contextualizes the results within the broader field of media bias detection and outlines potential avenues for system enhancement.

## 3.1 Project Outcomes

After the successful development and integration of the Bias Detection API and the Browser Extension, several performance evaluations and user experience tests were carried out. Below are the detailed outcomes observed in both components of the system.

## 1. BiasCheck API Results:

- **Response Time:** The API exhibited an average response time of 2–5 seconds for processing input text, depending on the length and complexity of the content. Although this response time is reasonable, further optimizations could be explored to reduce the latency.

- **Bias Classification:** The API classified bias labels (Left, Center, Right) based on the input text's political leanings, with an overall accuracy rate of ~91% when tested with a diverse set of real-world text samples.

- **Example Results:**

| Text Sample | Predicted Bias | Model Confidence (%) |
|---|---|---|
| "Government reforms hurt the marginalized the most." | Left | 87% |
| "The administration's efforts have been fair and effective." | Center | 92% |
| "Tax cuts empower the wealth creators of society." | Right | 89% |

Table 3.1 Sample Results of bias detection

## 2. Browser Extension Results:

The Chrome extension was developed to enhance user interaction by providing real-time bias detection on any webpage. Below are the important outcomes:

- **Real-Time Highlighting:** Upon visiting a page, the extension automatically scans visible text and highlights sections that display bias. This operation is seamless in most cases, though the time it takes to process the page varies. On average, it takes 2–5 seconds for the system to analyze and highlight the biased content.
- **Hover Information:** When users hover over a highlighted text, the extension provides the bias label (Left, Center, or Right), along with a confidence score, making it an interactive and informative experience. This ensures users understand why a specific piece of content has been flagged.

  While the extension's performance is quick, latency can be observed when processing larger pages, particularly those rich in content or when dynamic elements (e.g., infinite scroll) are present.

**Example Behavior:**

- ○ A biased sentence like "The government's policies are failing the poor" might get highlighted in red, indicating a Left-leaning bias. Hovering over it shows the text: "Bias: Left - Confidence: 88%."

- **User Feedback:** Testing across multiple news websites and blogs indicated positive user experiences with over 80% of users finding the system intuitive and informative. The highlighting behavior was seen as useful for promoting awareness about the potential biases in online content.

## 3.2 Discussion

The Bias Detection System demonstrated both functional effectiveness and room for improvement. Below are the key discussions based on the results observed:

### 1. Model Performance:

The BiasCheck-RoBERTa model exhibited strong performance in classifying political bias, with the following key observations:

- The model was highly accurate in classifying text into Left, Center, or Right categories, with 91% accuracy on the test set.

- Response Time: The 2–5 second delay for processing is reasonable for most applications, but as the system grows and more complex text is processed (e.g., longer documents), further optimization (such as using a faster model or introducing model quantization) might be necessary to reduce the delay.

### 2. Extension Functionality:

- **Real-Time Operation:** The extension operates well for detecting bias in real-time on most websites. However, there were slight delays (around 2–5 seconds) when processing pages with large amounts of text, dynamic content, or complex structures.

- **Usability:** Users reported that the hover feature was particularly helpful for understanding the reasoning behind the bias classification. It provided an engaging and insightful way to visualize bias detection results without disrupting the reading experience.

However, certain websites that use heavy JavaScript or infinite scrolling presented challenges in accurate text extraction, which could be mitigated by improving the extension's DOM parsing mechanism or enhancing it to work better with modern JavaScript-heavy sites.

## 3. Bias Detection Accuracy:

- Despite achieving an accuracy of ~91%, there are a few notable limitations:

  - **Complex Bias Interpretation:** Bias detection is inherently subjective, and while the model performs well on general political text, nuanced or subtle bias (e.g., sarcasm or highly contextual political references) was sometimes misclassified.

  - **Text Complexity:** The model tends to struggle with highly complex or highly opinionated texts. The hover mechanism and confidence score mitigate this, giving users a clear understanding of the model's decision-making process.

Figure 3.1 Evaluation Metrics

## 4. Performance in Different Scenarios:

- **News Articles:** The system performed well with news articles, as these tend to follow structured, formal language patterns, making it easier for the model to detect bias.

- **User-Generated Content (e.g., blogs, forums):** These types of content often contain more personal, varied opinions, and the model occasionally struggled to accurately classify bias in these informal formats.

## 3.3 Limitations and Future Directions

Despite promising results, there are several areas for improvement**:**

1. **Latency Optimization:** While the average processing time is acceptable for most users, reducing the 2–5 second delay would improve user experience, particularly for

those browsing long articles or news sites.

2. **Contextual Understanding:** Improving the model's ability to understand complex sarcasm, irony, or highly context-specific political commentary is crucial for increasing accuracy in real-world applications.

3. **Expansion of Supported Languages:** Currently, the system supports English text only. Future versions of the system could benefit from multilingual support to cater to a global audience.

4. **Extended Content Types:** The current system analyzes text-based content only. To enhance the functionality, future updates could involve the ability to analyze multimedia content (videos, images, etc.), which often play a crucial role in modern biased content.

# CHAPTER 4

# Conclusion & Future Enhancements

The Bias Detection System developed through this project offers a comprehensive framework to identify and highlight biases present in online content. The solution is composed of two tightly integrated components — a dedicated API and a browser extension — each playing a critical role in the system's functionality.

The API was designed as the analytical core of the project. It accepts text data as input, processes it through a series of Natural Language Processing (NLP) pipelines, and outputs a detailed bias score along with explanations. The API leverages state-of-the-art language models (such as fine-tuned versions of BERT, RoBERTa, or DistilBERT) combined with rule-based and statistical methods to detect linguistic patterns commonly associated with biased reporting — such as emotionally charged language, lack of source attribution, and slanted phrasing.

**Technologies used for the API:**

- Python (Flask/FastAPI) for building the web server

- NLP Libraries like spaCy, transformers (HuggingFace), and Scikit-learn

- Bias detection models trained on curated datasets (e.g., NewsBias, Wikipedia Neutrality Corpus)

- REST API architecture for easy integration

- Docker for containerized deployment

The browser extension acts as the client-facing layer. It interacts with the API by sending the textual content extracted from the currently viewed web page and displaying the analyzed results to the user in an unobtrusive but informative manner.

**Technologies used for the extension:**

- HTML/CSS/JavaScript for frontend development

- Manifest V3 to define the extension architecture

- Content Scripts to extract page data dynamically

- Background Scripts to handle API communication

- Browser Storage APIs to maintain user preferences

- Fetch/Axios libraries for sending data to the server

**Working Mechanism:**

1. The content script automatically scans and extracts the readable content from a loaded webpage.

2. The extracted text is cleaned and sent asynchronously via a POST request to the bias detection API.

3. The API processes the input, generates bias scores, flags key biased sentences, and returns the response.

4. The extension parses the API response and overlays bias indicators onto the original webpage.

5. The user can view detailed bias analysis by clicking on the extension popup.

Through careful integration of machine learning models, real-time communication protocols, and lightweight frontend design, the project successfully demonstrates the feasibility of

offering real-time bias detection during everyday browsing without significantly impacting performance or user experience.

Thus, this project contributes toward empowering users to develop critical reading skills and make informed judgments about the information they consume online.

**Future Enhancements**

Although the current system achieves its intended objectives, several improvements and extensions can be envisioned to elevate its effectiveness and robustness further:

**1. Model Upgradation and Continuous Learning**

- Incorporate newer transformer models such as RoBERTa or Longformer to better handle longer articles and subtle biases.

- Implement a continuous learning pipeline where user feedback on flagged content is used to retrain models, improving accuracy over time.

**2. Multilingual Support**

- Extend the API's capabilities to analyze non-English content, using multilingual models like XLM-RoBERTa.

- Build language detection pipelines to automatically route content through appropriate models.

**3. Deeper Contextual Analysis**

- Enhance bias detection by incorporating knowledge graphs and fact-checking APIs to detect not just linguistic bias but also factual inaccuracies.

- Use sentiment trajectory analysis to detect shifts in tone across an article.

**4. Explainable AI (XAI) Integration**

- Add features to explain why a particular piece of content was flagged as biased using techniques like LIME or SHAP.

- Provide sentence-level explanations to help users trust the system's judgments.

## 5. UI/UX Improvements

- Build a more interactive and user-friendly frontend with graphical bias indicators (e.g., heatmaps, dials, sliders).

- Allow user customization: for instance, setting personal thresholds for bias severity alerts.

## 6. Cross-Platform Extension Development

- Extend compatibility from Chrome to Firefox, Edge, and Safari.

- Release a mobile version of the browser extension for Android and iOS browsers.

## 7. Privacy and Security Enhancements

- Integrate end-to-end encryption for data sent between the browser extension and the API server.

- Implement local-first processing, minimizing data transmission to the server and preserving user privacy.

## 8. API Scalability and Optimization

- Shift from a simple REST API to a GraphQL API to reduce data overhead and offer more granular control over queries.

- Deploy the API on cloud-native architectures (like Kubernetes clusters) to achieve better load balancing and fault tolerance.

**9. Collaborative Bias Reporting**

- Allow users to submit bias reports collaboratively and crowdsource content evaluation.

- Build a community-driven bias index that rates sources based on collective user feedback.

**Final Thoughts**

This project lays the foundation for bias-aware browsing, a need of the hour in today's information-overloaded world. By combining AI, browser technologies, and human-centric design, it opens the pathway to critical consumption of digital content. Future enhancements, particularly in personalization, explainability, and scalability, can help turn this tool into a widely adopted standard for responsible internet usage.

# APPENDIX

## A.  RESEARCH PAPER

International Conference on Sustainable Computing and Communication Technologies : Submission (32) has been created.

External Inbox ×

Microsoft CMT <noreply@msr-cmt.org>
to me ▾

1:43 PM (0 minutes ago)

Hello,

The following submission has been created.

Track Name: ICSCCT2026

Paper ID: 32

Paper Title: Real-Time Political Bias Detection and Highlighting System for Web Content

Abstract:
This paper presents the design and deployment of BiasCheck, a real-time political bias detection system that analyzes website content, classifies it as left-leaning, centrist, or right-leaning, and highlights the text accordingly. The system uses a lightweight RoBERTa-based model fine-tuned for political bias classification, hosted on a FastAPI server, and integrated into a Chrome extension for seamless end-user experience. Experimental results demonstrate accurate classification and near-instantaneous feedback, making BiasCheck suitable for educational, media literacy, and journalism-related applications.

Keywords—Political Bias Detection, Chrome Extension, Natural Language Processing, FastAPI, Sentiment Analysis, Web Augmentation

Created on: Mon, 28 Apr 2025 08:12:58 GMT

Last Modified: Mon, 28 Apr 2025 08:12:58 GMT

Authors:
   – pk9598@srmist.edu.in (Primary)
   – ps3536@srmist.edu.in

Secondary Subject Areas: Not Entered

Submission Files:
   BiasCheck Research Paper.pdf (1 Mb, Mon, 28 Apr 2025 08:12:54 GMT)

Submission Questions Response: Not Entered

Thanks,
CMT team.

# B. SAMPLE CODING

```python
src >  train_no_preprocessing.py > ⊗ load_json_df
 2   import glob
 3   import json
 4   import pandas as pd
 5   import torch
 6   import numpy as np
 7   from datasets import Dataset
 8   from transformers import (
 9       RobertaTokenizer,
10       RobertaForSequenceClassification,
11       DataCollatorWithPadding,
12       Trainer,
13       TrainingArguments,
14   )
15   from sklearn.metrics import accuracy_score, f1_score
16
17   # ————— CONFIG ————————————————————————————————————————
18   # at the top, replace all your Windows-style paths with their WSL mounts:
19   DATA_DIR   = "/mnt/c/Users/peeka/Downloads/Politcial-Bias-Model/data"
20   JSON_DIR   = "/mnt/c/Users/peeka/Downloads/Politcial-Bias-Model/data/Article-Bias-Prediction-main/data/jsons"
21
22   # DATA_DIR = r"C:\Users\peeka\Downloads\Politcial-Bias-Model\data"
23   TXT_LABELS = {"Left Data": 0, "Center Data": 1, "Right Data": 2}
24   # JSON_DIR = r"C:\Users\peeka\Downloads\Politcial-Bias-Model\data\Article-Bias-Prediction-main\data\jsons"
25   MODEL_SAVE_PATH = "models/political-bias-model"
26   MODEL_NAME = "roberta-base"
27   # ————————————————————————————————————————————————————
28
29   def load_txt_df():
30       rows = []
31       for folder_name, label in TXT_LABELS.items():
32           # WSL path: .../data/Center Data/Center Data
33           folder = os.path.join(DATA_DIR, folder_name, folder_name)
34           if not os.path.isdir(folder):
35               print(f"⚠  Missing folder: {folder}")
36               continue
37           for fn in os.listdir(folder):
38               if fn.lower().endswith(".txt"):
39                   full = os.path.join(folder, fn)
40                   with open(full, encoding="utf-8") as f:
41                       txt = f.read().strip()
42                   if txt:
43                       rows.append({"text": txt, "label": label})
44       print(f"✅ Loaded {len(rows)} TXT samples")
45       return pd.DataFrame(rows)
46
47
48   def load_json_df():
49       rows = []
50       bias_map = {"left":0, "center":1, "right":2}
51       if not os.path.isdir(JSON_DIR):
```

```python
src >  train_no_preprocessing.py > ⊗ load_json_df
 63      def main():
 69          # 2) to HuggingFace Dataset and train/test split
 70          ds = Dataset.from_pandas(df[["text","label"]])
 71          ds = ds.train_test_split(test_size=0.2, seed=42)
 72          train_ds, eval_ds = ds["train"], ds["test"]
 73
 74          # 3) tokenizer & model
 75          tokenizer = RobertaTokenizer.from_pretrained(MODEL_NAME)
 76          model     = RobertaForSequenceClassification.from_pretrained(MODEL_NAME, num_labels=3)
 77
 78          # 4) tokenization
 79          def tok_fn(batch):
 80              return tokenizer(batch["text"], truncation=True, padding="max_length", max_length=512)
 81          train_ds = train_ds.map(tok_fn, batched=True)
 82          eval_ds  = eval_ds.map(tok_fn, batched=True)
 83
 84          # 5) data collator, metrics
 85          data_collator = DataCollatorWithPadding(tokenizer)
 86          def compute_metrics(p):
 87              preds = np.argmax(p.predictions, axis=1)
 88              return {
 89                  "accuracy": accuracy_score(p.label_ids, preds),
 90                  "f1": f1_score(p.label_ids, preds, average="weighted"),
 91              }
 92
 93          # 6) Trainer + training args
 94          device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
 95          print("Using device:", device)
 96          model.to(device)
 97
 98          training_args = TrainingArguments(
 99              output_dir="./results",
100              num_train_epochs=20,
101              per_device_train_batch_size=16,
102              per_device_eval_batch_size=16,
103              eval_strategy="epoch",
104              save_strategy="epoch",
105              learning_rate=2e-5,
106              weight_decay=0.01,
107              load_best_model_at_end=True,
```

```python
 63    def main():
 93        # 6) Trainer + training args
 94        device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
 95        print("Using device:", device)
 96        model.to(device)
 97
 98        training_args = TrainingArguments(
 99            output_dir="./results",
100            num_train_epochs=20,
101            per_device_train_batch_size=16,
102            per_device_eval_batch_size=16,
103            eval_strategy="epoch",
104            save_strategy="epoch",
105            learning_rate=2e-5,
106            weight_decay=0.01,
107            load_best_model_at_end=True,
108            save_total_limit=1,
109            push_to_hub=False,
110            logging_dir="./logs",
111        )
112
113        trainer = Trainer(
114            model=model,
115            args=training_args,
116            train_dataset=train_ds,
117            eval_dataset=eval_ds,
118            tokenizer=tokenizer,
119            data_collator=data_collator,
120            compute_metrics=compute_metrics,
121        )
122
123        # 7) train!
124        print("🖊 Starting training...")
125        trainer.train()
126        print("💾 Saving model to", MODEL_SAVE_PATH)
127        trainer.save_model(MODEL_SAVE_PATH)
128        tokenizer.save_pretrained(MODEL_SAVE_PATH)
129        print("✅ Training complete!")
130
131    if __name__ == "__main__":
132        main()
```

43

| SRM INSTITUTE SCIENCE AND TECHNOLOGY | | |
|---|---|---|
| **(Deemed to be University u/ s 3 of UGC Act, 1956)** | | |
| **Office of Controller of Examinations** | | |
| REPORT FOR PLAGIARISM CHECK ON THE DISSERTATION/PROJECT REPORTS FOR UG/PG PROGRAMMES **(To be attached in the dissertation/ project report)** | | |
| 1 | Name of the Candidate **(IN BLOCK LETTERS)** | **PRANAV KAUL** **PRANAY SHARMA** |
| 2 | Address of the Candidate | **Gautam Budh Nagar, Noida, Uttar Pradesh** **Bhartiya Nagar, Bilaspur, Chhattisgarh** |
| 3 | Registration Number | **RA2211028010078** **RA2211028010067** |
| 4 | Date of Birth | **05/04/2004** **03/10/2004** |
| 5 | Department | Department of Networking and Communications |
| 6 | Faculty | Engineering and Technology, School of Computing |
| 7 | Title of the Dissertation/Project | BiasCheck: A Political Bias Detector using RoBERTa |
| 8 | Whether the above project /dissertation is done by | Group **a)** If the project/ dissertation is done in group, then how many students together completed the project : 2 **b)** Mention the Name & Register number of other candidates : **PRANAV KAUL(RA2211028010078)** **PRANAY SHARMA(RA2211028010067)** |

| | | | | |
|---|---|---|---|---|
| 9 | Name and address of the Supervisor / Guide | **Mail ID: saravanm7@srmist.edu.in**<br><br>**Mobile Number: +91 89730 59349** | | |
| 10 | Name and address of Co-Supervisor / Co- Guide (if any) | **NIL** | | |

<br>

| | | | | |
|---|---|---|---|---|
| 11 | Software Used | Turnitin | | |
| 12 | Date of Verification | Apr 29, 2025, 2:58 PM GMT+5:30 | | |
| 13 | **Plagiarism Details: (to attach the final report from the software)** | | | |
| **Chapter** | **Title of the Chapter** | **Percentage of similarity index (including self citation)** | **Percentage of similarity index (Excluding self-citation)** | **% of plagiarism after excluding Quotes, Bibliography, etc.,** |
| **1** | Introduction | 12% | 12% | 11% |
| **2** | Literature Survey | 15% | 15% | 14% |
| **3** | Sprint Planning and Execution Methodology | 03% | 03% | 2% |
| **4** | Result and Discussion | 05% | 05% | 4% |
| **5** | Conclusion and Future Enhancement | 02% | 02% | 01% |
| **6** | | | | |
| **7** | | | | |
| **8** | | | | |
| **9** | | | | |
| **10** | | | | |
| **Appendices** | | | | |
| I / We declare that the above information have been verified and found true to the best of my / our knowledge. | | | | |

|  |  |
|---|---|
| **Signature of the Candidate** | **Name & Signature of the Staff (Who uses the plagiarism check software)** |
| **Name & Signature of the Supervisor/ Guide** | **Name & Signature of the Co-Supervisor/Co-Guide** |
| **Name & Signature of the HOD** | |

# Minor Project_report

## newclassificaition

- Quick Submit
- Quick Submit
- SRM Institute of Science & Technology

---

## Document Details

Submission ID
**trn:oid:::1:3232985202**

Submission Date
**Apr 29, 2025, 2:58 PM GMT+5:30**

Download Date
**Apr 29, 2025, 3:06 PM GMT+5:30**

File Name
**FINAL_PROJECT_REPORT_1.pdf**

File Size
**2.1 MB**

---

# 14% Overall Similarity

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

## Filtered from the Report

- Bibliography

## Match Groups

- **46** Not Cited or Quoted 13%
  Matches with neither in-text citation nor quotation marks
- **1** Missing Quotations 0%
  Matches that are still very similar to source material
- **2** Missing Citation 1%
  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

9% 🌐 Internet sources
4% 📖 Publications
12% 👤 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.

# REFERENCES

1. S. Raza, O. Bamgbose, V. Chatrath, S. Ghuge, Y. Sidyakin, and A. Y. Muaad, "Unlocking Bias Detection: Leveraging Transformer-Based Models for Content Analysis," arXiv preprint arXiv:2310.00347, 2023.

2. S. Suryawanshi, M. Ghosh, S. Chakraborty, and S. Ghosh, "A systematic review on media bias detection," Expert Systems with Applications, vol. 236, p. 121488, 2023.

3. A. Wessel, L. Leidner, and S. Schuster, "Improving Media Bias Detection with State-of-the-Art Transformers," in Proceedings of the 4th Workshop on NLP for Internet Freedom, 2022, pp. 1-11.

4. S. Raza et al., "Unlocking Bias Detection: Leveraging Transformer-Based Models for Content Analysis," in Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, 2023, pp. 1-12.

5. T. Wessel, L. Leidner, S. Schuster, and M. Reiche, "Experiments in News Bias Detection with Pre-Trained Neural Transformers," arXiv preprint arXiv:2406.09938, 2024.

6. N. Ji and Y. Sun, "Media Legitimacy Detection: A Data Science Approach To Locate Falsehoods And Bias Using Supervised Machine Learning And Natural-Language Processing," in Proc. 8th Int. Conf. Data Mining and Applications (DMAP), 2022.

7. P. Cáceres Ramos, J. M. Vázquez, V. P. Álvarez, and J. L. D. Olmedo, "I2C at PoliticEs 2022: Using Transformers to Identify Political Ideology in Spanish Tweets," in CEUR Workshop Proceedings, vol. 3202, 2022, pp. 49-58.

8. Minh Vu, "Political Bias Detection in News Articles using NLP and Deep Learning," Senior Thesis, Earlham College, 2018.

9. M. Sap, D. Card, S. Gabriel, Y. Choi, and N. A. Smith, "Social Bias Frames: Reasoning about Social and Demographic Bias in Language," in Proc. ACL, 2020, pp. 5477-5490.

10. A. Iyyer, P. Enns, J. Boyd-Graber, and P. Resnik, "Political Ideology Detection Using Recursive Neural Networks," in Proc. ACL, 2014, pp. 1113-1122.

11. K. Sim, A. Acre, and N. A. Smith, "Measuring Ideological Proportions in Political Speeches," in Proc. EMNLP, 2013, pp. 91-101.

12. J. Gangula, S. K. Sahu, and A. Kumar, "Headline Attention Network for Bias Detection in News Articles," in Proc. 16th Workshop on Asian Language Resources, 2019, pp. 67-76.

13. A. Soni, S. Soni, and S. K. Saha, "KnowBias: Domain Adaptation for Political Bias Detection," arXiv preprint arXiv:2010.05338, 2020.

14. Y. Liu et al., "RoBERTa: A Robustly Optimized BERT Pretraining Approach," arXiv preprint arXiv:1907.11692, 2019.

15. S. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in Proc. NAACL-HLT, 2019, pp. 4171-4186.

16. J. Pennington, R. Socher, and C. Manning, "GloVe: Global Vectors for Word Representation," in Proc. EMNLP, 2014, pp. 1532-1543.

17. T. Wolf et al., "Transformers: State-of-the-Art Natural Language Processing," in Proc. EMNLP: System Demonstrations, 2020, pp. 38-45.

18. J. Howard and S. Ruder, "Universal Language Model Fine-tuning for Text Classification," in Proc. ACL, 2018, pp. 328-339.

19. S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," Neural Computation, vol. 9, no. 8, pp. 1735-1780, 1997.

20. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," arXiv preprint arXiv:1810.04805, 2018