


Chapter 7

Highlighting the Challenges of Bias and Fairness in Hate Speech Detection

Manish Mittal

 <https://orcid.org/0000-0002-7029-9576>

Brainware University, India

Manish Tiwari

Sir Padampat Singhania University, India

Ruchi Doshi

ResAISHala Technocrats Pvt. Ltd., Udaipur, India

Kamal Kant Hiran

Sir Padampat Singhania University, India

ABSTRACT

This chapter offers a comprehensive overview of the challenges, advancements, and ethical concerns in automated hate speech detection. It critiques early keyword-based methods and highlights the shift toward advanced, context-aware, and culturally informed AI models. Key issues include bias in training data, subjectivity in annotation, and evolving adversarial tactics. The chapter emphasizes the essential role of human oversight and advocates hybrid human-AI moderation for balanced judgment. It stresses the importance of transparency, accountability, and fairness within complex legal and regulatory frameworks. Addressing technological limits such as multimodal data and annotation quality, it calls for interdisciplinary col-

DOI: 10.4018/979-8-3373-3063-1.ch007

Copyright © 2026, IGI Global Scientific Publishing. Copying or distributing in print or electronic forms without written permission of IGI Global Scientific Publishing is prohibited. Use of this chapter to train generative artificial intelligence (AI) technologies is expressly prohibited. The publisher reserves all rights to license its use for generative AI training and machine learning model development.

laboration and ongoing innovation. Ultimately, it argues that effective moderation must protect free expression and vulnerable groups, uniting technical strength with ethical responsibility.

INTRODUCTION

The problem of bias and fairness in the identification of hate speech is introduced in this chapter. At this point in time, the increasing dependence on automated programming to police content online makes it even more imperative to comprehend the processes by which computer systems are “trained” to identify and flag hate speech. The goal of this chapter is precisely to provide a clear scope of what hate speech is, and an analysis of why there has been so much bias and unfairness in the challenges of detecting it. Through the exploration of real-world twelve challenges in the development and deployment of hate speech detection technologies as well as real world examples of limitations, the use of secondary data, this chapter analyze uses both secondary data and a real-world example of the limitations and recommendation that twenty key challenges of developing and deploying hate speech detection technology.

The identification of hate speech on online platforms like Facebook, YouTube or Instagram poses a number of difficult challenges, issues and complexities with bias and fairness being the most salient among those. Among these challenges, one of the most central is the lack of a consensus definition of what constitutes “hate speech”. Determine what constitutes hate speech is culturally and legal specific, thus automated detection is problematic in multiple ways. In this chapter I will use a series of definitions and interpretations taken for granted within the scope of this study to attempt to understand some of the challenges surrounding the identification and regulation of hate content in digital environments (*The Digital Services Act Package | Shaping Europe’s Digital Future*, n.d.).

Hate speech encompasses any form of communication—be it spoken, written, or behavioural—that denigrates or discriminates against individuals or groups based on inherent characteristics such as race, ethnicity, religion, gender, sexual orientation, disability, or other identity factors. This includes the use of pejorative or discriminatory language aimed at a person or group based on who they are UNESCO (Countering hate speech, 2024). Key Characteristics of Hate Speech include forms of expression; hate speech can manifest through various mediums, including images, cartoons, memes, objects, gestures, and symbols, and can be disseminated both offline and online. Targets; it specifically targets individuals or groups based on attributes such as religion, ethnicity, nationality, race, colour, descent, gender, or other identity factors.

22 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/highlighting-the-challenges-of-bias-and-fairness-in-hate-speech-detection/393861

Related Content

Insider Attack Analysis in Building Effective Cyber Security for an Organization

Sunita Vikrant Dhavale (2019). *Cyber Law, Privacy, and Security: Concepts, Methodologies, Tools, and Applications* (pp. 1408-1425).

www.irma-international.org/chapter/insider-attack-analysis-in-building-effective-cyber-security-for-an-organization/228790

Privacy and Security: Safeguarding Personal Data in the AI Era

Geeta Sandeep Nadella, Hari Gonaygunta, Mohan Harishand Pawan Whig (2025). *Ethical Dimensions of AI Development* (pp. 157-174).

www.irma-international.org/chapter/privacy-and-security/359642

Hybrid Privacy Preservation Technique Using Neural Networks

R. VidyaBanuand N. Nagaveni (2019). *Cyber Law, Privacy, and Security: Concepts, Methodologies, Tools, and Applications* (pp. 542-561).

www.irma-international.org/chapter/hybrid-privacy-preservation-technique-using-neural-networks/228744

Developing Confidence Building Measures (CBMs) in Cyberspace Between Pakistan and India

Tughral Yamin (2019). *Cyber Law, Privacy, and Security: Concepts, Methodologies, Tools, and Applications* (pp. 1539-1602).

www.irma-international.org/chapter/developing-confidence-building-measures-cbms-in-cyberspace-between-pakistan-and-india/228798

Ethical Implications of the Techno-Social Dilemma in Contemporary Cyber-Security Phenomenon in Africa: Experience From Nigeria

Essien Essien (2019). *Cyber Law, Privacy, and Security: Concepts, Methodologies, Tools, and Applications* (pp. 1200-1213).

www.irma-international.org/chapter/ethical-implications-of-the-techno-social-dilemma-in-contemporary-cyber-security-phenomenon-in-africa/228778