

Chapter 12

Detecting Hate Speech in the Digital Age: AI Solutions for Real- Time Moderation

C. V. Suresh Babu

 <https://orcid.org/0000-0002-8474-2882>

Hindustan Institute of Technology and Science, India

S. Deva Darshini

 <https://orcid.org/0009-0001-7214-5175>

Hindustan Institute of Technology and Science, India

ABSTRACT

This chapter explores the application of artificial intelligence (AI) in real-time hate speech detection on social media platforms, aiming to address the limitations of traditional moderation techniques. Utilizing machine learning frameworks such as TensorFlow and PyTorch, the study examines the effectiveness of AI models in processing vast amounts of data to identify nuanced forms of hate speech. Key findings indicate that while AI can significantly enhance detection speed and accuracy, challenges remain regarding bias in training datasets and the need for human oversight to ensure ethical moderation practices. The implications of this research highlight the necessity for a balanced approach that integrates AI capabilities with human judgment to foster safer online environments while preserving free speech. This study contributes to ongoing discussions about the ethical deployment of AI technologies in content moderation.

DOI: 10.4018/979-8-3693-9904-0.ch012

1. INTRODUCTION

1.1 Overview of Hate Speech on Social Media

Social media apps like Facebook, Twitter, and YouTube have, over the past ten years, changed everything about how we communicate. But as all this was wrought upon the digital sphere, harming behaviors increased with hate speech catching one of the highest places among the most prevalent. Hate speech that is mostly confined in cyberspace addresses particular individuals based on the different dimensions of race, religion, gender, or sexual orientation with psychological damage, social division, and even material harm. Social media's anonymity and mass following establish an easy venue for the diffusion of toxic messages, which can transform possibly private thoughts to viral content in just a minute. This has caught the attention of social media providers and governments alike, raising very important questions about how such platforms should properly be managed and moderated. That kind of traditional human moderation has proved woefully inadequate given the number of posts that must be read daily to be able to adequately moderate (Davidson, T., Warmlesley, D., Macy, M., & Weber, I. 2017).

1.2 The Importance of Real-Time Detection and Response

Hence, real-time hate speech detection is the need of the hour. Compared to the static contents, social media posts are highly dynamic in nature. The malicious tweet or Facebook post may get retweeted, shared, and read by millions within the shortest span of time. The posts may cause irreparable damage to people and communities if not detected and moderated in due time. Real-time detection will allow hate speech to be identified right after it is posted, hence such widespread cannot occur. Reactions can also forestall the proliferation of the content pretty fast, so some of the potential damage is limited. Apart from acceleration of moderation, a real-time system is also a form of deterrence. Users may think twice about posting harmful content if they know that it can be removed within a couple of seconds. Without such ability, social media organizations would merely react and play catch-up to hate speech once done (Schmidt, A., & Wiegand, M. 2017).

1.3 Role of AI in Hate Speech Detection

One of the most powerful tools that emerged was through artificial intelligence because the problem is indeed massive for real-time hate speech moderation; in classical systems based on simple keyword matches or user reports, AI models can process large volumes of data, recognize patterns, and, to some extent, understand

30 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/detecting-hate-speech-in-the-digital-age/371741

Related Content

Capturing the Context of Concepts using the Transaction Graph through a Mobile NHS Case Study

Ivan Launders (2016). *International Journal of Conceptual Structures and Smart Applications* (pp. 35-47).

www.irma-international.org/article/capturing-the-context-of-concepts-using-the-transaction-graph-through-a-mobile-nhs-case-study/171390

Design and Usage of a Process-Centric Collaboration Methodology for Virtual Organizations in Hybrid Environments

Thorsten J. Dollmann, Peter Loos, Michael Fellmann, Oliver Thomas, Andreas Hoheisel, Peter Katranuschkov and Raimar Scherer (2011). *International Journal of Intelligent Information Technologies* (pp. 45-64).

www.irma-international.org/article/design-usage-process-centric-collaboration/50485

Efficient Identification of Structural Relationships for XML Queries using Secure Labeling Schemes

S. Sankari and S. Bose (2016). *International Journal of Intelligent Information Technologies* (pp. 63-80).

www.irma-international.org/article/efficient-identification-of-structural-relationships-for-xml-queries-using-secure-labeling-schemes/171441

A Multimodal Sentiment Analysis Model for Graphic Texts Based on Deep Feature Interaction Networks

Wanjun Chang and Dongfang Zhang (2024). *International Journal of Ambient Computing and Intelligence* (pp. 1-19).

www.irma-international.org/article/a-multimodal-sentiment-analysis-model-for-graphic-texts-based-on-deep-feature-interaction-networks/355192

Leveraging Artificial Intelligence for Advanced Threat Detection and Autonomous Response in Cybersecurity

Venkata Ramana Kaneti, D. Jayasutha, B. Yamini, M. G. Dinesh, L. A. Anto Gracious and P. Girija (2026). *Implementing Enterprise Cybersecurity With AI* (pp. 135-160).

www.irma-international.org/chapter/leveraging-artificial-intelligence-for-advanced-threat-detection-and-autonomous-response-in-cybersecurity/395163