# Facilitating Anonymous Communication on Social Networks via AI-Driven Content Moderation

Siddhesh Mengade[1*], Pranjali Chopade[1], Parth Tate[1], Shraddha Patil[1]

[1]Computer Engineering, Genba Sopanrao Moze College of Engineering

pentosid@gmail.com, chopadepranjali3@gmail.com, parthtate575@gmail.com, shraddhajpatil93@gmail.com

## Abstract

*The Anonymous Messaging Application (Anonify) is designed to provide a secure platform for users to send and receive anonymous feedback, messages, or questions, fostering open dialogue in academic, professional, and social settings. The primary objective of this project is to address the need for a safe space where individuals can express their thoughts without fear of judgment or retribution, while ensuring robust content moderation to prevent misuse. To achieve this, Anonify leverages modern technologies such as Next.js for seamless server-side rendering, Natural Language Processing (NLP) for real-time content analysis, Auth.js for secure user authentication, and Tailwind CSS for an intuitive user interface. The methods employed include the integration of an AI-powered moderation engine that uses machine learning algorithms to detect and filter inappropriate content, such as abusive language, spam, or malicious intent. NLP models are utilized to identify patterns associated with bullying or fraudulent behavior, while fraud detection mechanisms monitor for repeated misuse attempts. These measures ensure that the platform maintains a respectful and constructive environment for users.*

*The results demonstrate Anonify's ability to effectively moderate content in real-time, filtering out harmful messages while allowing constructive feedback to reach the intended recipients. The application successfully balances anonymity with accountability, providing a secure space for open communication without compromising user safety.*

*Thus Anonify represents a significant step forward in anonymous communication platforms, addressing the challenges of unregulated online environments. By integrating advanced AI moderation and user-friendly design, the application promotes constructive dialogue, enhances personal and professional relationships, and ensures a safe, anonymous feedback system.*

**Keywords:** Anonymous Communication, Schema, Validation, Auth.js, Zod, Spam & Fraud Detection, Feedback Collection, Content Moderation, Next.js, Foul Language Detection.

## 1. INTRODUCTION

Digital communication has revolutionized how we interact, yet it often limits honest dialogue due to fears of judgment and exposure. In many environments—such as workplaces, educational settings, and personal networks—the absence of secure, anonymous channels discourages individuals from sharing candid opinions. This gap not only stifles innovation and creative

problem-solving but also hampers personal and professional growth by promoting a culture of silence around sensitive topics.

To address these challenges, we propose Anonify, an anonymous messaging platform designed to empower users to express their thoughts freely without the fear of identity disclosure. Anonify leverages advanced, mobile-first technology alongside an AI-powered content moderation system that filters out inappropriate, offensive, or spam messages in real time. This combination of anonymity and real-time safeguarding is intended to create a secure space that promotes open and constructive dialogue.

The primary objectives of Anonify are threefold:

- Enhance User Communication: Provide a secure platform for anonymous messaging that encourages the exchange of genuine feedback.
- Promote Constructive Feedback: Foster an environment where candid, respectful, and meaningful conversations can thrive.
- Integrate AI-Powered Moderation: Ensure user safety and maintain a respectful digital space by automatically filtering harmful content.

We hypothesize that by removing the barriers associated with identity disclosure, Anonify will not only increase the quality and frequency of feedback but also stimulate innovative thinking and personal development. In essence, our approach posits that secure anonymity, when coupled with advanced moderation, can transform digital interactions and create a more open, supportive communication landscape.

This paper is organized into nine sections. Section 1 introduces the concept of Anonify and its significance in fostering anonymous communication along with the problem statement, the objectives, underlying hypotheses, and motivation driving the development of Anonify. Section 2 provides a literature survey, reviewing existing research and platforms related to anonymous communication. Section 3 details the system architecture and design of Anonify, including its database structure and AI-powered moderation system. Section 4 discusses the progress in implementation, focusing on backend development and API routing and also explores the use of AI for suggesting messages and content moderation. Section 5 evaluates the results and impact of the platform. Finally, Section 6 concludes the paper with a summary of findings and future work.

## 2. LITERATURE SURVEY

The concept of anonymous feedback platforms has evolved significantly over the years, with early platforms like Ask.fm and NGL.link paving the way for modern solutions. These platforms were designed to facilitate open communication by allowing users to share feedback anonymously.[1] However, they faced significant challenges, particularly related to cyberbullying and misuse, which highlighted the need for improved security measures and user-friendly interfaces. Research on anonymous feedback platforms emphasizes the importance of robust content moderation and

secure authentication mechanisms to ensure a safe and constructive environment for users. This is where Anonify differentiates itself by integrating advanced technologies such as Zod for real-time database management, Auth.js for seamless and secure user authentication, and AI algorithms for sentiment analysis, spam detection, and fraud prevention.[2][3] These features enable Anonify to moderate anonymous feedback effectively, ensuring a safe and respectful communication environment.

**Evolution of Anonymous Communication Platforms:** Anonymous feedback platforms have been studied extensively for their impact on social behavior, with research indicating both positive and negative outcomes. While anonymous messaging can promote open communication and honest feedback, it also carries risks such as cyberbullying and online harassment. For instance, platforms like NGL.link and Sarahah gained popularity for their ability to facilitate anonymous communication but faced criticism for their inability to effectively address misuse.

- NGL.link (Not Gonna Lie): Inspired by earlier anonymous platforms, NGL.link gained traction by integrating with Instagram, allowing users to receive anonymous feedback via stories. While it offered moderation tools to filter inappropriate content, privacy concerns were primarily addressed on the receiver's side. However, the platform introduced a subscription model that allowed receivers to view the identity of senders, effectively undermining the concept of anonymity.
- Sarahah: This platform was designed to encourage constructive feedback but became associated with online harassment, leading to its eventual removal from app stores. The backlash against Sarahah underscored the importance of content moderation and user safety in anonymous communication platforms.

The following research papers were reviewed to gain insights into anonymous communication and its associated concepts.

## 2.1    WhisperLink: A Novel Anonymous Messaging Service for a Secured Data Communication (2024 IEEE/ACIS)

A recent advancement in anonymous communication is WhisperLink, a messaging service designed to enhance privacy in digital communication. Developed by Morales et al., WhisperLink operates on Google Cloud and allows users to create temporary chat rooms that self-destruct after 24 hours, ensuring confidentiality without requiring logins. The platform employs end-to-end encryption to ensure that only intended recipients can read messages and uses security questions for additional authentication. By not storing data beyond 24 hours, WhisperLink guarantees private and transient conversations, making it suitable for social, professional, and private interactions.[4]

**Key Features of WhisperLink:**

- End-to-End Encryption: Ensures that only the intended recipients can read messages, preventing unauthorized access.

- No-Login Requirement: Users can communicate without providing personal information or creating accounts, enhancing privacy.
- Temporary Chat Rooms: Messages and chat rooms automatically delete after 24 hours, ensuring no residual data is stored.

**Relevance to Anonify**: The no-login requirement implemented in WhisperLink has inspired a similar feature in Anonify, particularly for anonymous senders. This approach minimizes the amount of personal data stored by the platform, reducing the risk of identity exposure and data leaks.

## 2.2    The Impact of Anonymity on Communication in the Metaverse (2024 IEEE COMPSAC)

Another significant study by Nakayama and Sumi explores the impact of anonymity on communication in the metaverse. The researchers used virtual avatars to examine how varying levels of anonymity affect interactions. Four distinct communication groups were tested:

- Anonymous-Anonymous: Both participants communicated without revealing their identities.
- Anonymous-Non-Anonymous: One participant remained anonymous, while the other was identifiable.
- Non-Anonymous-Anonymous: One participant was identifiable, while the other remained anonymous.
- Non-Anonymous-Non-Anonymous: Both participants communicated without anonymity.

The study found no significant differences in immersive communication across the groups, suggesting that anonymity did not influence communication behavior. However, the researchers concluded that visual anonymity plays a critical role in preventing rude or violent behaviors in virtual environments. This finding implies that even superficial anonymity, where personal details are not shared, can foster respectful communication in immersive settings like the metaverse.[5]

**Relevance to Anonify**: The study's findings are particularly relevant to Anonify, as the platform operates in a virtual environment where negative behaviors like rudeness or violence could be prominent. To address this, Anonify incorporates AI algorithms to detect and filter harmful content in real-time. The anonymous-to-non-anonymous communication scheme explored in the study has also influenced the design of Anonify, ensuring a balanced approach to anonymity and accountability.

These studies collectively highlight the advancements, challenges, and opportunities in anonymous communication, providing a strong foundation for the development of secure and user-centric platforms like Anonify.

**Key Takeaways from Literature Survey:** The literature survey highlights several critical insights:

- Content Moderation: Robust content moderation is essential to prevent misuse and ensure a safe environment for users.
- Privacy and Security: Features like end-to-end encryption and no-login requirements enhance user privacy and reduce the risk of data leaks.
- Balancing Anonymity and Accountability: While anonymity promotes open communication, it must be balanced with mechanisms to prevent abuse, such as AI-powered moderation and secure authentication.
- Adaptability to New Environments: The findings from studies on the metaverse underscore the importance of adapting anonymous communication systems to new digital environments, ensuring they remain effective and secure.

By integrating these insights, Anonify aims to create a platform that not only facilitates anonymous communication but also ensures user safety and constructive interactions.

**Changes Made in Reformatted Version:**

- Logical Flow: The section now flows logically, starting with an overview of anonymous feedback platforms, followed by a discussion of specific platforms (NGL.link, Sarahah), and then diving into recent research (WhisperLink, Metaverse study).
- Contextual Connections: Each subsection connects to the broader theme of anonymous communication and explains how the findings are relevant to Anonify.
- Clear Headings: Subheadings (e.g., "Evolution of Anonymous Feedback Platforms," "WhisperLink: A Novel Anonymous Messaging Service") make the section easier to navigate.
- Relevance to Anonify: Each study or platform discussed is explicitly tied back to how it informs or influences the design and functionality of Anonify.

# 3. SYSTEM ARCHITECTURE & DESIGN DETAILS

The lack of anonymous communication channels hinders open dialogue, stifles innovation, and limits personal growth, as individuals often fear expressing honest opinions on sensitive topics. Without anonymity, important discussions are left unaddressed, and creative thinking is suppressed. To solve this, the anonymous messaging receiver application provides a secure platform for users to share thoughts freely while ensuring privacy. By integrating AI-powered content moderation, it filters inappropriate or harmful content, promoting a respectful environment where users can engage in honest, meaningful conversations without fear of judgment or privacy concerns.
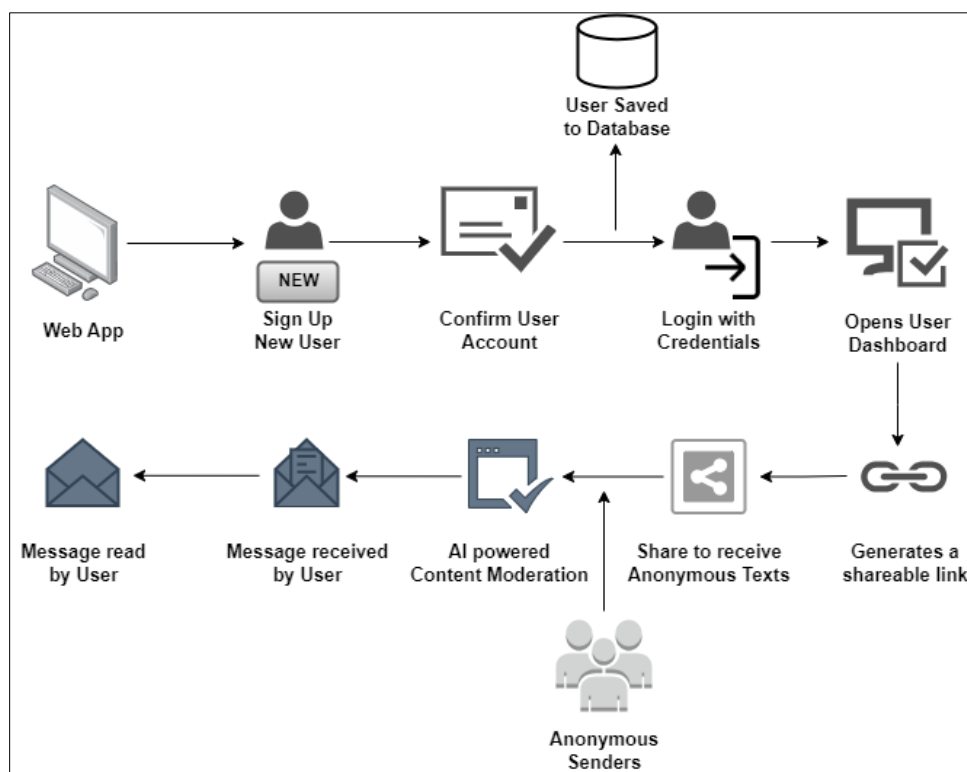
## 3.1 System Architecture



*Figure 1 System Architecture*

The diagram above represents the system architecture demonstrating the flow of an anonymous messaging application. Here's a brief information about the working process:

- User Registration: A new user signs up, they confirm their account via a verification code sent on their e-mail. Once verified, their data is saved in the database.
- Login & Dashboard: Users log in and access their dashboard, where they generate a shareable link which can be used to receive messages from anonymous senders.
- Link Sharing: The user shares the link, allowing anonymous senders to submit messages.
- AI Content Moderation: Messages undergo AI-powered moderation to filter inappropriate content in real time. Thus, blocking any inappropriate content and not sending it to the user at all.
- Message Reception: The user receives, views, and reads messages (that pass the AI Content Moderator) through the platform.

## 3.2 Database Design

The two schemas present in the database are, User Table and Message Table. And they outline the database schema for the anonymous messaging web application. Following explanation describes these schemas.

6

The two schemas present in the database are, User Table and Message Table. And they outline the database schema for the anonymous messaging web application. Following explanation describes these schemas.

The User Table manages the user-related data such as login credentials, account status (verified or not), and whether the user is open to receiving anonymous messages. Whereas the Message Table manages the messages that are sent anonymously to users. Each message is linked to a user through a foreign key (user_id), and the content is stored securely.
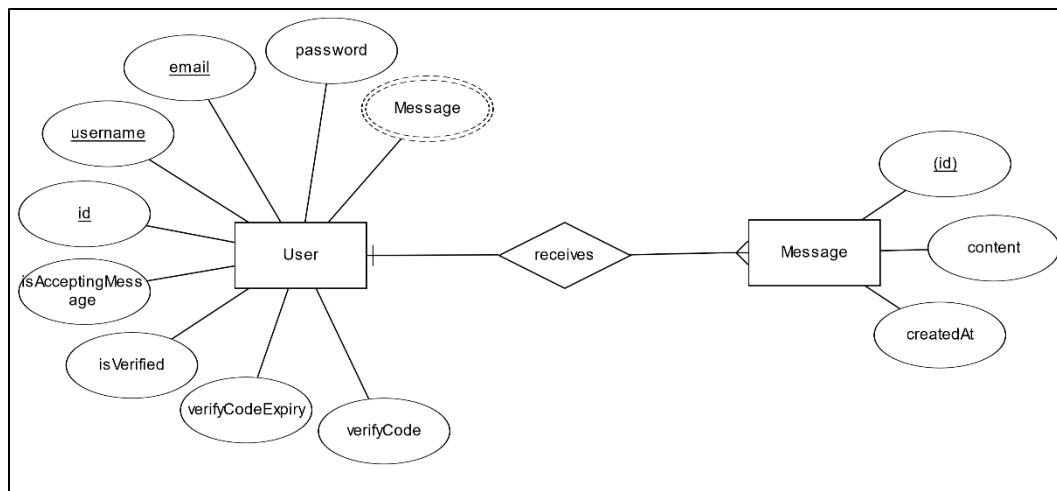


*Figure 2 ER Diagram*

## 3.3   AI-Powered Content Moderation

In the future, AI can play a crucial role in moderating content sent anonymously to non-anonymous receivers on the platform. A key method involves utilizing Sentiment Analysis to detect and filter harmful or inappropriate messages before they reach the recipient. AI models can analyze the sentiment of each message and classify it as positive, neutral, or negative, enabling the system to automatically flag offensive, abusive, or inappropriate messages.

**Key Benefits of AI-Based Moderation:**

- Real-time Detection: AI algorithms like Naive Bayes, SVM, or LSTM can analyze messages instantly, flagging harmful content in real time.
- Automated Filtering: AI uses supervised learning to automate moderation, flagging extreme negative content like hate speech for blocking or review.
- Context Awareness: NLP models like BERT and GPT understand message context, reducing false positives by distinguishing between harmful and harmless expressions.
- Scalability: Deep learning models and RNNs enable efficient processing of large message volumes, ensuring timely moderation even with high traffic.
- Adaptation to Evolving Language: Word embeddings (e.g., Word2Vec, FastText) allow AI to adapt to new slang and evolving language, ensuring continued accuracy in detecting harmful content.

**Additional AI Techniques for Moderation:**

- Anomaly Detection: Algorithms like K-Means clustering, Isolation Forest, or Autoencoders can be used to detect unusual patterns in messaging behavior, such as a sudden increase in negative content or signs of bullying, and flag these as potential threats.
- Contextual Relevance: AI models like BERT or T5 (Text-to-Text Transfer Transformer) can assess whether a message aligns with the overall conversation, helping to identify messages that are out of context or disruptive, such as spam or malicious content.[6]

**Limitations of AI-Powered Content Moderation**

- Evolving Language & Slang: Language evolves continuously, introducing new slang and offensive expressions. AI models trained on static datasets struggle to recognize emerging harmful phrases and adapt dynamically.[7]
- Over Reliance on Lexicons: Lexicon-based models depend on predefined word lists, often failing to grasp context. They may incorrectly flag innocent words or miss offensive variations across languages.[8]

By incorporating AI-based content moderation techniques powered by advanced algorithms such as transformers, RNNs, SVMs, and deep learning models, the platform can ensure that messages sent anonymously are carefully filtered for harmful content before being delivered to non-anonymous receivers. This will create a safer and more positive user experience, reducing the likelihood of inappropriate or abusive interactions.

# 4. PROGRESS IN THE IMPLEMENTATION OF BACKEND FUNCTIONALITY

## 4.1 Backend Development

In Next.js, API routing and route handlers are fundamental for building server-side functionality.

### 1. API Routes:

API Routes allow you to create backend endpoints in your Next.js application, all within the /src/api directory. Each file in this directory becomes an API endpoint. These routes enable handling various HTTP methods such as GET, POST, PUT, DELETE, etc. Example: A file at src/api/sign-up.ts becomes available at the endpoint /api/sign-up.

- Request & Response: Next.js API routes provide access to the req and res objects from Node.js, allowing you to handle requests and responses easily.
- Middleware: You can add middleware to authenticate users or protect routes.
- Error Handling: You can handle errors by returning appropriate HTTP status codes.
- Limitations: API routes don't support features like request body parsing by default; you need to manually parse them for certain content types.

## 2. Route Handlers:

Route Handlers in Next.js are specific to the new app directory architecture, allowing more granular control over how a route responds to different HTTP methods. These handlers can be used to define custom logic for different HTTP verbs, such as GET, POST, and PUT, within the same file. They allow for easy handling of dynamic routes by using parameters. Example: You can define route handlers by creating a file in app/api/route.js, which handles HTTP requests.[9]

- Streaming and Buffering: Route handlers support both, making it easier to handle large amounts of data, such as file uploads or real-time updates.
- Advanced Capabilities: They also offer advanced features like incremental static regeneration, server-side rendering, and dynamic rendering based on request parameters.

In summary, API routes provide a simple way to handle server-side logic in a Next.js app, while route handlers offer more granular control over handling HTTP requests in the newer app directory structure.

## 4.2 SUGGESTING MESSAGES USING AI

GPT-3 and GPT-Neo are both transformer-based language models designed for natural language processing tasks. They are built on the transformer architecture, which uses self-attention mechanisms to understand relationships within a sequence of words.

- GPT-3: Developed by OpenAI, GPT-3 is a state-of-the-art model with 175 billion parameters. It excels in generating human-like text across a variety of tasks due to its extensive training on diverse datasets. GPT-3 employs autoregressive decoding, where each word or token is generated sequentially, considering the input and previously generated tokens. This model is pre-trained on large datasets and can be fine-tuned for specific applications, such as text summarization or creative writing.[10]
- GPT-Neo: An open-source alternative to GPT-3, GPT-Neo, developed by EleutherAI, is designed to make large-scale language modeling accessible to the community. While not as powerful as GPT-3 in terms of parameter size, GPT-Neo models like the 2.7 billion parameter variant perform well on tasks like text generation and question answering. They are trained on public datasets and can be integrated for custom applications.[11]

## 4.3 Algorithms Used for Suggesting Messages:

- Transformer Architecture: Both models rely on self-attention mechanisms, which allow them to capture long-range dependencies in text efficiently. This makes the models adept at understanding the context of a conversation or a user prompt.
- Autoregressive Decoding: In message suggestion tasks, the models generate tokens sequentially. The generation of each token depends on the previously generated tokens and the input prompt. This ensures the output is coherent and contextually appropriate.

9

- Prompt Conditioning: The models are conditioned with a well-crafted prompt (e.g., "Generate three open-ended questions"). This guides the generation process to align the output with the desired structure and intent, such as creating engaging social media messages.
- Tokenization: Before processing, input text is divided into smaller units called tokens. The model processes these tokens to predict the next token iteratively, forming a complete response.
- Beam Search or Sampling: For diverse and creative outputs, methods like Top-k Sampling or Nucleus Sampling are often employed. These introduce randomness by selecting from the top-ranked tokens at each step, ensuring variability in generated messages while maintaining relevance.

## 4.4 CONTENT MODERATION USING AI

The GPT-3.5 Turbo model, which is available through platforms like Hugging Face, is primarily based on a transformer architecture, specifically the GPT (Generative Pre-trained Transformer) model, which uses unsupervised learning techniques to understand and generate human-like text. While GPT-3.5 Turbo itself is not explicitly designed for offensive content detection, it can be adapted for that task using various AI techniques.[12]

- Transformer Architecture (GPT-based models)
  - Self-Attention Mechanism: The core component of transformer models like GPT is the self-attention mechanism, which allows the model to consider the entire context of the text while making predictions. This context-awareness is useful when detecting subtle offensive content that might not be obvious in isolated words but is evident in the overall tone or context.
  - Pre-training and Fine-tuning: GPT-3.5 is pre-trained on a massive amount of text data and fine-tuned on specific tasks like offensive language detection. During fine-tuning, the model is trained on datasets labeled with offensive and non-offensive content, allowing it to learn patterns and features specific to harmful language.
  - Binary Classification: This is the most common technique where the model is trained to classify text as either "offensive" or "non-offensive." The training data consists of labeled examples, where each text instance is tagged with one of these two labels. Popular algorithms for text classification include:
    - Logistic Regression
    - Naive Bayes
    - Support Vector Machines (SVM)
    - Neural Networks (especially deep learning models like CNNs and RNNs)
  - Multi-class Classification: Instead of just two categories, the model can classify text into multiple categories such as "hate speech," "abusive," "profane," or

10

"non-offensive." This allows for more nuanced detection of different types of offensive content.

- Text Classification
  - Supervised Learning: In the case of offensive language detection, supervised learning techniques are often employed. The model is trained on a labeled dataset of text examples, where each example is tagged as "offensive" or "non-offensive." The model learns to classify text based on these labels.
  - Multi-Class Classification: Sometimes, the offensive content detection model might categorize text into multiple classes, such as "abusive," "profane," "hate speech," or "non-offensive." These models use softmax or similar activation functions to output probabilities for each category.
- Sentiment Analysis
  - While sentiment analysis typically categorizes text as positive, negative, or neutral, it can be adapted to detect offensive language by associating specific negative sentiments or aggressive tones with harmful content. The model might be fine-tuned on datasets where aggressive or hateful speech is labeled as "negative sentiment," triggering the detection of offensive language.
  - Lexicon-Based Approaches: In these approaches, a predefined dictionary (or lexicon) of words associated with specific sentiments (e.g., positive, negative, neutral) is used. Sentiment scores are calculated based on the presence of these words in the text. Examples include:
    - VADER (Valence Aware Dictionary and sEntiment Reasoner): It is a lexicon and rule-based sentiment analysis tool specifically tailored for social media text. VADER detects not just positive or negative sentiment but also the intensity and polarity of emotions, which can be useful in detecting offensive or aggressive content.[13]
    - SentiWordNet: A lexical resource for sentiment analysis that assigns sentiment scores to words. This can be used to identify negative emotions or offensive language.[14]
  - Rule-Based Sentiment Analysis: Rules are created to identify sentiment based on patterns in text. For example, the presence of certain words (e.g., "hate," "violence," "abuse") might trigger a "negative" sentiment label. This approach can be useful in detecting offensive or hateful language.

# 5. CONCLUSION

In conclusion, the development of Anonify represents a significant step toward creating a secure and user-friendly anonymous messaging platform. While we have made substantial progress in implementing core functionalities, including user registration, message sending, and basic content moderation, there is still work to be done to enhance the application further. Our focus on integrating AI-powered content moderation has laid the groundwork for a safer user experience, effectively filtering out inappropriate content and ensuring a positive environment for our users.

And that is what we aim to work on for the next stage of development. Moreover, the project has provided valuable insights into the complexities of full-stack development, particularly in leveraging Next.js for seamless server-side rendering and efficient API routing. Although some advanced features, such as OTP verification, are yet to be fully integrated, our progress thus far underscore the project's potential.

Hence, Anonify strives to create a platform that not only enables anonymous communication but also prioritizes user safety and trust in digital interactions. By striking a balance between anonymity and accountability, Anonify provides a secure environment where users can freely express themselves while fostering positivity. As the platform continues to evolve, we are committed to enhancing its features and ensuring it serves as a catalyst for open, safe, and meaningful communication.

# 6. REFERENCES

1. Y. Lee and K. H. Kim, "De-motivating employees' negative communication behaviors on anonymous social media: The role of public relations," Public Relations Review, vol. 46, no. 4, p. 101902, 2020, doi: 10.1016/j.pubrev.2020.101902.
2. T. Milosevic et al., "Effectiveness of artificial intelligence–based cyberbullying interventions from youth perspective," Social Media + Society, 2023, doi: 10.1177/20563051231156469.
3. J. Pratt, "Hugely Popular NGL App Offers Teenagers Anonymity In Comments About Each Other," Forbes, Jun. 29, 2022. [Online]. Available: https://www.forbes.com/sites/iainmartin/2022/06/29/hugely-popular-ngl-app-offers-teenagers-anonymity-in-comments-about-each-other/. [Accessed: Feb. 11, 2025].
4. M. Morales, A. Boyina, D. Kothari, M. Moniruzzaman, and A. Sultana, "WhisperLink: A novel anonymous messaging service for a secured data communication," in 2024 IEEE/ACIS International Conference on Computer and Information Science (ICIS), 2024.
5. Y. Nakayama and K. Sumi, "The impact of anonymity on communication in the metaverse," in Proceedings of the IEEE COMPSAC, 2024.
6. J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in Proc. NAACL-HLT 2019, Minneapolis, MN, USA, Jun. 2019, pp. 4171–4186.
7. P. Fortuna and S. Nunes, "A Survey on Automatic Detection of Hate Speech in Text," in ACM Computing Surveys (CSUR), vol. 51, no. 4, pp. 1–30, Jul. 2018.
8. T. Davidson, D. Warmsley, M. Macy, and I. Weber, "Automated Hate Speech Detection and the Problem of Offensive Language," in Proceedings of the International AAAI Conference on Web and Social Media (ICWSM), vol. 11, pp. 512–515, Jun. 2017.
9. R. Gorwa, "Content moderation, algorithmic governance, and the new politics of the internet," Internet Policy Review, vol. 8, no. 2, 2019, doi: 10.14763/2019.2.1414.
10. T. B. Brown, B. Mann, N. Ryder, M. Subbiah, et al., "Language Models are Few-Shot Learners," in Advances in Neural Information Processing Systems, vol. 33, pp. 1877–1901, 2020.

11. EleutherAI, "GPT-Neo," GitHub repository. [Online]. Available: https://github.com/EleutherAI/gpt-neo. [Accessed: Feb. 11, 2025].

12. B. Liu, Sentiment Analysis and Opinion Mining. Morgan & Claypool Publishers, 2012.

13. C. J. Hutto and E. Gilbert, "VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text," in Proc. Int. AAAI Conf. Web Soc. Media (ICWSM), 2014, pp. 216–225.

14. A. Esuli and F. Sebastiani, "SentiWordNet: A Publicly Available Lexical Resource for Opinion Mining," in Proc. 5th Int. Conf. Lang. Resour. Eval. (LREC), Genoa, Italy, 2006, pp. 417–422.