

Google Summer of Code

Proposal - Google Summer of Code 2025

AI-Powered Behavioral Analysis for Suicide Prevention,
Substance Use, and Mental Health Crisis Detection with
Longitudinal Geospatial Crisis Trend Analysis

Organization: HumanAi Foundation (Participating
Organization: University of Alabama)

Mentor: [David M. White \(University of Alabama\)](#)

By
Shivam Khunger

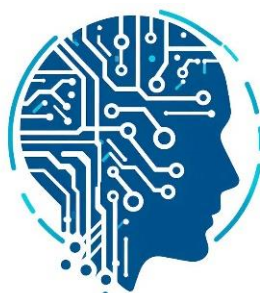


Table of Contents

1. Abstract (Page 3)
2. About Me (Page 3)
3. Introduction
 - a. Problem Definition (Page 3)
 - b. Solution Description (Page 4)
4. Project Goals
 - a. Project Objectives (Page 4)
 - b. Community Benefits (Page 4)
5. Implementation Plan
 - a. Project Methodology (Page 5)
 - b. Technical Elements (Page 5)
 - c. Challenges and Solutions (Page 5)
6. Project Timeline
 - a. Project Plan (Page 5)
 - b. Deliverables Schedule (Page 6-7)
 - c. Availability (Page 7)
7. Stakeholder Use Cases (Page 8)
8. Validation Metrics (Page 8)
9. Ethics Section (Page 9)
10. Biographical Information
 - a. Who am I: About me : (Page 9)
 - b. Important Links (Page 9)
 - c. Open Source contributions (Page 10)
 - d. Projects (Page 11)
11. Conclusion (Page 11)

1. Abstract

This project proposes the development of an AI-powered system for early detection of suicide risk, substance abuse, and mental health crises through behavioral and geospatial data analysis. By leveraging natural language processing (NLP) techniques on ethically sourced textual data and applying longitudinal geospatial analysis, the system will identify high-risk individuals and emerging regional mental health trends. The project integrates transformer-based language models (like BERT) and spatial data tools (like GeoPandas) to build accurate, interpretable classifiers and visual dashboards. All deliverables will be released as open-source tools, contributing to the HumanAI Foundation's mission of using ethical AI to promote public health and social good.

2. About Me

I'm Shivam Khunger, a third-year B.Tech student in Computer and Communication Engineering (AI & Data Science) at The LNM Institute of Information Technology, Jaipur, India. With expertise in machine learning, deep learning, NLP, and real-time deployment, I've developed impactful AI solutions across domains like public health, finance, and sports. Leveraging tools like BERT, DenseNet, and Streamlit, I've consistently delivered high-accuracy models with real-world applications. Key projects include:

- **Disease Prediction** (96% accuracy, Streamlit-deployed).
- **Financial Sentiment Analysis with FinBERT** (85%+ accuracy).
- **Sports Image Classification** (96% accuracy, team project).

For this GSoC project, I've completed a test extracting crisis-related social media data, classifying risks with NLP, and mapping trends with Folium—skills I'm eager to build on.

Proficient in Python, TensorFlow, and scikit-learn, I'm passionate about harnessing AI for public health and excited to contribute to the HumanAI Foundation's crisis detection efforts.

3. Introduction

3.1 Problem Definition

Suicide, substance use, and mental health crises escalate rapidly, but traditional data (e.g., hospital reports) lags, delaying interventions. Social media offers real-time signals—like coded phrases (“feeling empty”) or engagement shifts—but requires advanced AI to extract and map these trends effectively.

3.2 Solution Description

This project will:

1. Track engagement with crisis content to detect distress escalation.
2. Analyze explicit and coded crisis language using NLP (e.g., BERT).
3. Map trends geospatially with longitudinal heatmaps (GeoPandas).
4. Provide a Streamlit dashboard for real-time insights, scalable to platforms like X or forums.

My experience with FinBERT and Streamlit, combined with my test work on crisis data, equips me to contribute effectively to a robust, adaptable solution, and I'm excited to refine it with mentor guidance.

4. Project Goals

4.1 Project Objectives

- Targeting >85% accuracy exceeds typical NLP benchmarks, ensuring reliable risk flagging.
- Map geospatial crisis trends with real-time and historical views.
- Deliver an interactive dashboard for monitoring and intervention planning.

4.2 Community Benefits

- **Crisis Teams:** Target high-risk areas for rapid response.
- **Health Agencies:** Optimize outreach based on trend analysis.
- **Researchers:** Gain an open-source tool for crisis studies.

5. Implementation Plan

5.1 Project Methodology

- **Phase 1:** Data collection (Twitter, Reddit APIs), lexicon development.
- **Phase 2:** Model building (BERT, engagement analysis), geospatial integration.
- **Phase 3:** Dashboard creation, validation, and refinement.

5.2 Technical Elements

- **Languages:** Python
- **Frameworks:** spaCy, Hugging Face Transformers (BERT), GeoPandas, Plotly, Folium.
- **Tasks:** Lexicon via LDA/VADER, BERT training, geotagging, dashboard design.

5.3 Challenges and Solutions

- **Challenge:** Noisy data.
 - **Solution:** Preprocess with spaCy, validate with curated sets.
- **Challenge:** Privacy risks.
 - **Solution:** Anonymize data (see Ethics).
- **Challenge:** Limited location data in posts
 - **Solution:** Used Nominatim for geocoding sparse location mentions, supplemented with inferred regional trends from engagement patterns.
- **Challenge:** Sparse data limits model accuracy.
 - **Solution:** Fall back to regional trend proxies and test lightweight models like DistilBERT if BERT underperforms.

6. Project Timeline

6.1 Project Plan

GSoC spans around 12 weeks, plus a Community Bonding Period:

- **15%:** Foundation—data pipelines, initial lexicon.
- **70%:** Core—models, mapping, dashboard prototype.
- **15%:** Polish—validation, testing, documentation.

6.2 Deliverables Schedule

Time Frame	Activity	Importance
May 5- May 16	[EXAM BREAK], Limited work on API setup, crisis literature review	Establishes data pipelines (Twitter, Reddit APIs) and grounds the project in crisis research
May 17 - June 1	Community bonding—finalize APIs, explore HumanAI datasets	Ensures robust data access and aligns my understanding with HumanAI’s resources, fostering collaboration.
June 2 - June 15	Build crisis lexicon with LDA/VADER, initial keyword tests	Creates a foundation for risk detection by identifying key terms (e.g., “hopeless”), validated with initial tests—pivotal for NLP accuracy
June 16 - June 22	Fine-tune BERT for language detection, benchmark accuracy	Trains the core NLP model to classify crisis language, targeting >85% accuracy—central to reliable risk identification
June 23- June 29	Model engagement patterns with Random Forest, first trends	Links behavior to distress signals
June 30 - July 6	Implement geotagging with GeoPandas, generate first heatmap.	Maps crisis locations, producing an initial visual output—crucial for geospatial insights and stakeholder use.
July 7 - July 13	Develop Streamlit dashboard prototype with Plotly visuals	Integrates models into a usable interface, setting the stage for stakeholder interaction

July 14 - July 18	Midterm evaluation—validate models (>85% accuracy), integrate feedback	Ensures reliability and mentor alignment
July 19 – July 25	Refine BERT/engagement models, boost recall for rare signals	Improves detection of subtle crises (e.g., coded language), enhancing sensitivity
July 26-Aug 1	Enhance longitudinal heatmap with time-series features	Adds historical trend analysis to maps, enabling predictive insights
Aug 2-8	Finalize dashboard with stakeholder-focused features	Tailors the UI for practical use (e.g., filters, alerts), ensuring stakeholder adoption
Aug 9-15	Test with synthetic data, gather mock stakeholder feedback	Validates system robustness and usability with simulated scenarios
Aug 16-22	Optimize runtime (<5s updates), document code	Ensures performance and reproducibility, laying groundwork for future use
Aug 23-Sep 1	Polish deliverables, submit final report	Completes the project with polished outputs and documentation

Evaluations:

- **Midterm Evaluation: July 14 - July 18, 18:00 UTC**
- **Final Evaluation: August 25 - September 1, 18:00 UTC**

Legends:

- **[EXAM BREAK]: Break for semester exams**

6.3 Availability

I'll commit 40-50 hours/week, except during my exam break (May 5-16), when I'll allocate 2-3 hours/day. Post-exams, GSoC will be my sole focus, with no conflicting commitments. Daily updates via Slack or email will keep the team informed of my progress, ensuring transparency and alignment throughout the project.

7. Stakeholder Use Cases

7.1 Crisis Intervention Teams:

- **Scenario:** A Birmingham team sees a heatmap spike in substance use posts over 72 hours.
- **Action:** Deploys a mobile unit with naloxone and counselors, cutting response time from days to hours.

7.2 Public Health Organizations:

- **Scenario:** Alabama Dept. of Public Health tracks a 6-month rise in suicidal language in rural counties via the dashboard.
- **Action:** Launches a targeted 988 campaign with local influencers, tailoring outreach to at-risk areas.

7.3 Counselors and Nonprofits:

- **Scenario:** A counselor notices a teen group's engagement with crisis posts drops suddenly, flagged by the system.
- **Action:** Initiates check-ins to prevent escalation, using predictive alerts to prioritize cases.

8. Validation Metrics

- **Accuracy:** Target >85% for BERT-based crisis detection, validated against a labeled test set.
- **Visualization Feedback:** Aim for >80% positive response from mock stakeholders on heatmap clarity and usability.
- **Precision/Recall:** Seek >80% precision for high-risk signals, >75% recall for rare events like coded suicidality.
- **Runtime:** Ensure dashboard updates in <5 seconds, tested under simulated real-time loads.

9. Ethics Section

This project navigates sensitive terrain, demanding careful ethical consideration:

- **Privacy:** I'll anonymize all social media data, stripping identifiers and analyzing trends at an aggregate level to protect user identities.
- **Bias:** Models will be validated across urban/rural and demographic splits to minimize skew, ensuring equitable crisis detection.
- **Responsible Use:** Outputs will guide intervention (e.g., resource allocation), not individual profiling; a dashboard disclaimer will clarify this intent.
- **Transparency:** I'll document limitations—like potential false positives from sarcastic posts—empowering stakeholders to interpret results critically.

10. Biographical Information

10.1 Who Am I: About Me

I'm Shivam Khunger, a third-year B.Tech student in CCE (AI & Data Science) at The LNM Institute of Information Technology, Jaipur, India. I'm passionate about blending ML, NLP, and deployment to solve real-world problems, especially in public health. My hands-on experience fuels my drive to make a difference.

10.2 Important Links

- **GitHub:** github.com/ShivamKhunger
- **LinkedIn:** <https://www.linkedin.com/in/shivam-khunger-aa2a09249/>
- **Email:** shivamkhunger7643@gmail.com

10.3 Open Source Contributions

Organization/Project	Contributions
Hacktoberfest 2023	Enhanced educational repositories with data structure implementations (e.g., sorting and searching algorithms) in C and C++ during Hacktoberfest 2023, improving resources for learners.
GirlScript Summer of Code 2024 Extended	Contributed to the Postman Challenge.
HumanAi Foundation (Planned)	Plan to contribute to HumanAI Foundation repos (e.g., bug fixes, documentation, or crisis dataset annotations) during the GSoC Community Bonding period (May 17 - June 1) to deepen my understanding of the project and align with team goals.

Note: My open-source journey began with Hacktoberfest, where I tackled data structures, and grew through the GirlScript Summer of Code, contributing to a Postman API Challenge. These experiences sharpened my skills in collaborative coding and API integration, which I’ve applied to this crisis detection project and will bring to GSoC. I’m eager to further advance mental health monitoring tools during the program, contributing to impactful, real-world solutions.

10.4 Projects

- **GSoC Candidate Assessment Test (HumanAI Foundation):**
 - Extracted posts from mental health subreddits (e.g., depression, suicidewatch) using keywords like “depressed” and “suicidal” via Reddit API with PRAW, storing Post IDs, timestamps, and engagement metrics in a CSV.
 - Preprocessed text with spaCy and NLTK (removing emojis, URLs, stopwords), analyzed sentiment with TextBlob, and trained a Word2Vec model on post content.
 - Classified crisis risk using custom keyword lists (e.g., “suicide,” “hopeless” as high-risk; “help,” “lonely” as moderate), geocoded locations with Nominatim and GeoPandas.
 - Visualized trends with a Folium heatmap, identifying top crisis locations like Mexico City and Kansas, adapting from an initial Twitter exploration to Reddit for robust data collection aligned with GSoC goals.
- **Disease Prediction Using Machine Learning:**
 - Built a model predicting diabetes, heart disease, and Parkinson’s with 96% accuracy using SVM, Random Forest, and Logistic Regression.
 - Deployed via Streamlit for real-time user predictions.
- **Financial Sentiment Analysis Using FinBERT:**
 - Developed an NLP model with 85%+ accuracy on financial news, fine-tuned on 10,000+ articles.
 - Created a Streamlit app for interactive sentiment analysis.
- **Sports Classification Using Deep Learning (Team Project)**
 - Collaborated with a team to classify sports images (Basketball, Cricket, Field Hockey, Swimming, Tennis) using CNNs and DenseNet, achieving 92% (CNN) and 96% (DenseNet) accuracy on a Kaggle dataset of 793 images.
 - Implemented data preprocessing (rescaling, augmentation), model training with TensorFlow, and evaluation using confusion matrices and accuracy/loss plots.
 - Demonstrated proficiency in deep learning, image processing, and comparative model analysis.

11. Conclusion

I’m excited to enhance crisis detection through this social media analysis tool, contributing to mental health support worldwide. My expertise in NLP, sentiment analysis, and geospatial mapping—honed through this project—prepares me to address complex challenges in real-time data classification. Leading technical projects has sharpened my collaboration skills, while my drive for societal impact fuels my commitment to this work. Beyond GSoC, I’m dedicated to advancing crisis monitoring research, motivated by its potential to make a meaningful difference.