# <PU CODE HACKATHON 3.0 2026>

| Fields | Information |
| --- | --- |
| **Problem Statement Title** | AI Powered Online Harassment Detector |
| **Team Name** | Cipher Trace |
| **Team Leader Name** | Jyanesh Naidu |
| **Institute Name** | Parul Institute of Engineering and Technology |
| **Track Name** | Cyber Security |
| **Team Member** | Manish Kashyap, Susritha Swamyvari ,Koushik Yadav |

## Proposed Solution
## AI-Powered Online Harassment Detection System

| | |
|---|---|
| ❌ **Problem** | ▪ Rapid rise in online harassment and cyberbullying<br>▪ Manual moderation is slow and not scalable<br>▪ Harmful impact on users' mental health<br>▪ Lack of real-time detection systems |
| 🎯 **Objective** | ▪ Detect online harassment using AI in real time<br>▪ Classify abuse types and severity<br>▪ Prevent repeated harassment incidents<br>▪ Ensure safer digital platforms |
| ⚙️ **Approach** | ▪ Analyze user messages using NLP techniques<br>▪ Apply AI models for harassment classification<br>▪ Generate risk score for severity detection<br>▪ Take automated actions (warn, block, report) |

| USE CASES | TECHNOLOGY STACK | |
|---|---|---|
| 👥 Social Media Platforms | 🤖 AI / Machine Learning | ▪ Python<br>▪ DistilBERT (NLP Model)<br>▪ HuggingFace Transformers<br>▪ Scikit-learn |
| 🎓 Online Education Platforms | ⚙️ Backend Services | ▪ FastAPI<br>▪ REST APIs<br>▪ JWT Authentication |
| 🎮 Gaming & Live Streaming Platforms | 🖥️ Frontend & Dashboard | ▪ React.js / Streamlit<br>▪ HTML, CSS, JavaScript |
| 🏢 Corporate Communication Tools | 🗄️ Database & Storage | ▪ MongoDB<br>▪ JSON-based logs |
| 🧠 Mental Health & Child Safety | 🔒 Security & Monitoring | ▪ Role-based access control<br>▪ Secure API communication |
| | ☁️ Deployment | ▪ Render / Railway<br>▪ Docker |

# PROCESS FLOWCHART

```
Start → User Sends Message / Comment → Text Preprocessing
                                        Cleaning, Tokenization,
                                        Stopword Removal         → AI / NLP Model Analysis
                                                                   DistilBERT Classification
```

```
Risk Score Calculation ← Harassment Type Detection ←
0–100 Severity Level      Hate Speech / Bullying /
                          Threat / Safe
```

```
                        Decision Engine
```

| Safe | Medium Risk | High Risk | Repeat Offender |
|------|-------------|-----------|-----------------|
| Allow Message | Issue Warning to User | Block & Report Content | Auto Ban User |

```
                    Admin Dashboard Update
                              ↓
                             End
```

# DEPENDENCIES

- ❑ **Limited availability of high-quality labeled harassment datasets**
- ❑ **Difficulty handling multilingual content, slang, and evolving abusive language**
- ❑ **Challenges in detecting sarcasm and context-based harassment**
- ❑ **High computational requirements for real-time moderation**
- ❑ **Risk of false positives affecting genuine users**
- ❑ **Need for continuous model updates to track new abuse patterns**

## Mitigation Approach

- ➢ **Use pre-trained transformer models**
- ➢ **Continuous learning using feedback loop**
- ➢ **Multilingual model support**
- ➢ **Human-in-the-loop moderation for edge cases**