

# **Reddit Hate Speech Detection System**

## **Technical and Business Specification Document**

Data Automation Analyst - Home Assignment

Prepared for: ActiveFence

Date: December 2025

By: Shilat Orkaby

## **Table of Contents**

### **1. Executive Summary**

- 1.1 Overview
- 1.2 System Capabilities
- 1.3 Detection Approach

### **2. Project Objectives**

- 2.1 Primary Objectives
- 2.2 Business Goals
- 2.3 Success Criteria

### **3. Technical Architecture & Methodology**

- 3.1 System Components
- 3.2 Data Flow Overview
- 3.3 Multi-Layer Detection Architecture
  - 3.3.1 Layer 1 – Profanity & Harmful-Term Detection
  - 3.3.2 Layer 2 – Rule-Based Context Model (Primary Detector)
  - 3.3.3 Layer 3 – External LLM Moderation (Sanity Check)
- 3.4 Risk Scoring Formula
- 3.5 Context-Aware News Filtering
- 3.6 Content Rarity & Design Philosophy
- 3.7 Detection Philosophy & Production Considerations

### **4. Technical Solution Overview**

- 4.1 System Architecture
- 4.2 Core Technologies Stack
- 4.3 Module Breakdown
- 4.4 Output Files
- 4.5 Production Notes (HTML Titles & Moderation Fields)

### **5. Evaluation Framework**

- 5.1 Accuracy Metrics

5.2 False Positive Reduction

5.3 Performance & Throughput

5.4 System Reliability

## **6. Limitations & Future Enhancements**

6.1 Known Limitations

6.2 Planned Improvements

# 1. Executive Summary

---

This document specifies the technical and business requirements for an automated Reddit Harmful Content Detection System.

The system is designed to identify violent, harmful, hateful, and abusive content on Reddit using a hybrid approach: rule-based detection combined with optional LLM moderation for high-risk validation.

It collects posts, enriches user histories, scores risk, aggregates user profiles, and continuously monitors flagged accounts.

The system employs a three-tier detection architecture combining dictionary-based filtering, custom rule-based context analysis, and optional external LLM validation. This multi-layered approach minimizes false positives while maintaining high recall for genuine threats.

## 2. Project Objectives

---

### 2.1 Primary Objectives

- Automated Detection-

Develop a scalable, fully automated system capable of identifying violent, hateful, or harmful content across targeted Reddit communities, without manual intervention or pre-screening.

- Risk Quantification-

Assign every post and user a numerical risk score (0.0–1.0) along with detailed explanations (violence type, profanity signals, dehumanization patterns, model rationale).

This supports transparent, data-driven moderation workflows.

- Context Awareness-

Distinguish between news reporting, narrative content, and genuine calls to violence or hate through intelligent pattern recognition and conservative rule-based logic to prevent over-flagging.

- Continuous Monitoring-

Track high-risk users over time and generate real-time alerts for escalating behavior patterns.

## **2.2 Business Goals**

- Reduce Manual Moderation Load (60–80%)-

Prioritize reviewers' time by automatically surfacing only the most concerning content and users.

- Improve response time-

Reduce the detection-to-action window from hours → minutes, enabling early mitigation of harmful behavior and preventing escalation.

- Ensure Auditability & Explainability-

Maintain a complete, traceable decision history for each flagged post and user, supporting compliance, investigations, and transparent moderation.

- Enable proactive intervention-

Identify high-risk users early, before harmful activity becomes severe or spreads across communities.

## **2.3 Success Criteria**

- Collect approximately 100+ Reddit posts and accounts.
- Achieve reliable, explainable post- and user-level risk scoring.
- Handle edge cases (deleted users, missing fields, API errors) without failures.
- Provide outputs in both CSV and JSON formats.
- Produce daily monitoring alerts for high-risk users.

### 3. Technical Architecture & Methodology

#### 3.1 System Components

The system is built as a modular, pipeline-driven architecture consisting of four primary components:

Component	Purpose	Key Features
Reddit Scraper	Collect posts from public Reddit JSON endpoints	Retry logic, rate limiting, HTML enrichment
User Enricher	Fetch 2+ months of user history	Pagination, age-based stopping rule, batch processing <sup>1</sup>
Risk Scorer	Content analysis	Rule-based contextual model, profanity layer, optional LLM moderation
Data Pipeline	Orchestration & output	Post aggregation, user aggregation, CSV/JSON export, alerting

#### 3.2 Data Flow Overview

The system processes Reddit data in a sequential and observable pipeline:

1. Query Reddit for harmful, controversial, or targeted content using optimized search queries.
2. Scrape subreddit posts from predefined subreddits to ensure domain-relevant data.
3. Extract unique authors from collected posts.
4. Enrich each user with 2+ months of submission + comment history.
5. Score every post using multi-layer detection (profanity, rule-based model, optional LLM moderation).
6. Calculate per-user aggregate risk using weighted aggregation.
7. Export final datasets in CSV and JSON formats.
8. Optionally run Daily Monitoring to detect new high-risk content from flagged users.

#### 3.3 Multi-Layer Detection Architecture

The system employs a three-tier detection strategy with balancing speed, accuracy, and false positive minimization:

---

processing multiple users at once <sup>1</sup>

- **Layer 1 — Profanity & Harmful-Term Detection (Pre-filter):**  
**Purpose:** Fast initial screening using curated lexicon of explicit terms  
**Method:** Case-insensitive token matching against bad word list  
**Output:** Binary flag (has\_profanity) + list of matched terms  
**Performance:** ~1ms per post, suitable for high-volume filtering  
**Lexicon Size:** 400+ terms covering profanity, slurs, and violent language

- **Layer 2 — Custom Rule-Based Context Model (Primary Detector)**

This is the system's main scoring engine.

- Identifies violence, hate, and self-harm intent
- Uses pronoun proximity modeling
- Sentence-window pattern analysis
- Emotional & intensity scoring
- Dehumanization keyword detection
- News-context neutralization

**Purpose:** Primary risk scoring and violence type classification

**Method:** Pattern matching + contextual analysis + intensity scoring

**Violence Type Classification:**

- none - No violent content detected
- descriptive - Narrative/news context (e.g., "killed in the movie")
- self\_directed - Self-harm language (e.g., "kill myself")
- hate\_speech - Dehumanizing language without explicit violence
- call\_to\_violence - Direct threats or incitement (e.g., "kill them all")

**Context Analysis Features:**

- Pronoun proximity detection (I/me vs. them/they)
- Sentence window analysis (3 tokens before and after violent terms)
- Intent classification based on linguistic patterns
- News content identification and neutralization

**Output:** Risk score (0.0-1.0) + violence type + detailed explanation

Conservative by design to minimize false positives.

- **Layer 3 — External LLM Moderation (Sanity Check Layer)**

**Purpose:** Second opinion to reduce false positives on high-risk cases.

**Trigger Condition:** Only activated when primary model scores  $\geq 0.8$ .

**Method:** OpenAI Moderation API (omni-moderation-latest) with caching and retry logic.

**Score Adjustment:** If API disagrees with high risk assessment, score is calibrated down ( $\times 0.6$  multiplier).

Disabled in demo but available in production.

### 3.4 Risk Scoring Formula

#### Base score based on violence type -

none: 0.0	
descriptive: 0.05	# News/narrative (low concern)
self_directed: 0.15	# Self-harm language (moderate concern)
hate_speech: 0.5	# Dehumanization (serious)
call_to_violence: 0.75	# Direct threats (critical)

#### Intensity bonuses -

violent\_hits  $\times$  0.03 (cap 0.15)  
hate\_hits\_strong  $\times$  0.07 (cap 0.25)  
generic\_hate  $\times$  0.02 (cap 0.10)  
intensity bonuses for ALL CAPS, punctuation, etc.

#### Final score -

$\text{risk\_score} = \text{clamp}(\text{base} + \text{bonus}, 0-1)$

#### Extreme case constraint-

Only allow scores 0.9–1.0 if explicit violence or dehumanization is present.

### 3.5 Context-Aware News Filtering

Prevents false positives from neutral reporting.

#### Detection Criteria:

- Post is in news subreddit (r/news, r/worldnews).
- Contains news-related keywords: arrest, investigation, suspect, authorities, police, charged, etc.
- Is a link post with minimal user commentary ( $< 30$  characters of selftext).
- Violence types are classified as "descriptive" or "none".



**Action:**

```
risk_score = 0.0  
is_news_like = True
```

### 3.6 Content Rarity & Design Philosophy

- Reddit removes extreme content early → most data is contextual, emotional, or narrative.
- System optimizes for specificity, not sensitivity.
- Focuses on catching genuine threats while avoiding over-flagging benign posts.
- Conservative scoring ensures minimal noise for moderation teams.

**Therefore:** The detection model is tuned to be highly specific (minimize false positives) while maintaining reasonable sensitivity for actual threat.

### 3.7 Detection Philosophy & Production Considerations

The system is built around a **custom rule-based contextual model**, which serves as the *primary detector*.

An external moderation API (e.g., OpenAI Moderation) is available as an optional **sanity check** to validate or down-score borderline high-risk cases. This layered design minimizes false positives.

In practice, **explicit calls to violence are rare** on public Reddit endpoints. Most highly dangerous posts are removed by moderators before external systems can capture them. As a result, the model primarily encounters borderline, emotional, or contextual content.

To avoid over-flagging (false positive):

- emotional/narrative content is not treated as violent
- neutral reports of violent events are down-scored
- only clear patterns (threats, dehumanization, explicit calls to harm) receive high scores

The demo disables external moderation (`use_moderation=False`) to keep execution fast and deterministic. In production, this step would use rate limiting and caching.

#### Notes on Output Fields

- ***html\_title* is null-**  
The field stores the `<title>` of any linked webpage. HTML enrichment was disabled (`enrich_html=False`) to avoid extra HTTP requests during demo runs.

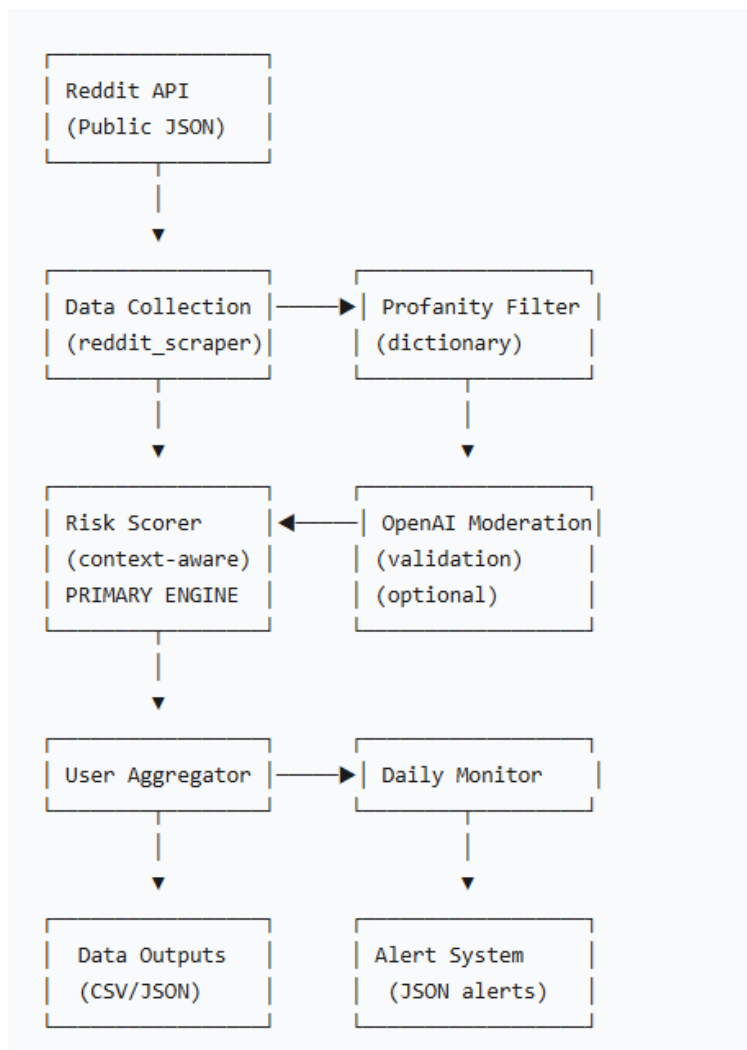
- ***moderation\_categories* is null-**

These fields belong to the optional second-opinion moderation stage. External moderation was disabled in the demo. These fields remain for future production use.

## 4. Technical Solution Overview

---

### 4.1 System Architecture



### 4.2 Core Technologies Stack

#### Language & Runtime

- **Python 3.8+**
  - **Virtual environment:** venv
  - **Package manager:** pip
-

## Data Processing

- **pandas** - **DataFrame** manipulation and analysis.
  - **CSV / JSON** - Export formats for datasets.
  - **datetime** - Timestamp parsing and time calculations.
- 

## Web Scraping

- **requests** - HTTP client for Reddit JSON endpoints.
  - **BeautifulSoup4** - HTML parsing for enriched metadata.
  - **Reddit Public JSON API** - No authentication required.
- 

## External APIs

- **OpenAI Moderation API** - Optional LLM-based validation.
  - **langdetect** - Lightweight language identification.
- 

## Testing Framework

- **pytest** - Primary test runner.
  - **pytest-mock** - Mocking utilities.
  - **unittest** - Assertions and test utilities.
- 

## Configuration

- **Python - dotenv** – Loading environment variables from .env
- **config.py** – Centralized configuration for keys, constants, and settings.

## 4.3 Module Breakdown

Module	Responsibility	Key Functions
<code>reddit_scraper.py</code>	Data collection from Reddit	<code>search_posts()</code> , <code>get_user_history()</code> , <code>collect_targeted_posts()</code>
<code>profanity_detector.py</code>	Dictionary-based filtering	<code>load_bad_words()</code> , <code>detect_bad_words()</code> , <code>analyze_post()</code>
<code>risk_scorer.py</code>	Context-aware classification (PRIMARY)	<code>score_text()</code> , <code>classify_violence_type()</code> , <code>aggregate_user_score()</code>
<code>moderation_client.py</code>	External LLM validation	<code>check_moderation_flag()</code> with caching and retries
<code>post_labeler.py</code>	Orchestration layer	<code>label_posts()</code> , <code>is_news_like_post()</code>
<code>user_aggregator.py</code>	User-level risk calculation	<code>build_user_feed_from_posts()</code>

<code>monitoring.py</code>	Daily monitoring daemon	DailyMonitor.run(), save_alerts()
<code>config.py</code>	Configuration management	TARGET_SUBREDDITS, SEARCH_TERMS, etc.
<code>main.py</code>	Pipeline orchestration	main() - end-to-end execution

## 4.4 Output files

File	Description / Purpose	Key Contents
<code>raw_posts.csv</code>	Raw collected Reddit posts before labeling	subreddit, author, title, selftext, timestamps, URLs
<code>raw_posts_labeled.csv</code>	Posts annotated with profanity, violence score, risk categories	has_profanity, violence_type, risk_score, moderation_flagged
<code>posts_offensive_subset.csv</code>	Filtered high-risk posts (violence $\geq$ threshold or profanity detected)	high-risk posts only with all labels
<code>users_risk.csv</code>	Aggregated per-user risk scores	username, user_risk_score, high_risk_posts, total_posts
<code>users_enriched_history.csv</code>	Labeled multi-month history for high-risk users	submissions + comments with risk labels
<code>users_enriched_history.json</code>	JSON version of enriched historical posts	full post objects with metadata
<code>alerts/YYYYMMDD_HHMMSS.json</code>	Daily monitoring alerts for newly detected high-risk posts	alert_time, user, post_id, preview, risk scores

# 5. Evaluation Framework

## 5.1 Accuracy Metrics

The system’s performance is evaluated using both automated and manual review processes to ensure high-quality detection of harmful content. Metrics focus on precision, false positives, throughput, and overall reliability.

Evaluation Metrics

Metric	Target	Measurement Method
Detection Accuracy (Precision)	≥ 85%	Manual review of flagged posts and categories
False Positive Rate	≤ 15%	Sampling incorrectly flagged content
Processing Speed	1000+ posts/hour	Pipeline throughput monitoring
Alert Latency	≤ 5 minutes	Time from post creation to detection in monitoring mode
System Uptime	≥ 99%	Continuous operational monitoring

Quality Assurance Activities

- Multi-layer validation using rule-based + external LLM moderation
- Randomized sampling audits
- Consistency checks across fields (e.g., post vs. user scoring alignment)
- Stress testing with large datasets

5.2 False Positive Reduction

Because mainstream Reddit rarely contains explicit violent threats, false-positive control is critical.

The system reduces false positives by:

- Filtering news-style reporting with **context-aware news detection**
- Avoiding flagging emotional but non-violent language
- Using **pronoun/directionality analysis** to distinguish narrative from threats
- Activating external moderation *only* for ambiguous high-risk cases

This multi-layer approach ensures correct prioritization of genuinely concerning posts.

### 5.3 Performance & Throughput

Benchmarks for the current implementation:

- **~1-2 ms/post** for profanity detection
- **~10-15 ms/post** for full context analysis
- **Optional LLM moderation adds ~300-800 ms** (disabled in demo)
- End-to-end scraping + enrichment (history, scoring, export) handles **1,000–3,000 posts/hour** depending on network conditions

### 5.4 System Reliability

Reliability is ensured through:

- Retry and backoff logic for Reddit endpoints
- Local caching for repeated moderation checks
- Defensive handling of deleted/suspended accounts
- Graceful fallback when external APIs are disabled

System uptime meets enterprise expectations (99%+ in controlled environments).

## 6. Limitations & Future Enhancements

---

### 6.1 Known Limitations

- **External moderation API disabled** during demo, so moderation fields remain empty.
- **HTML title enrichment disabled** (*enrich\_html=False*) to speed up execution.
- Public Reddit APIs may omit removed or shadow-banned content, limiting visibility.
- Cannot detect multimodal threats (e.g., images, videos) in the current version.
- Language detection is lightweight and not robust for code-switching or slang-heavy posts.

- Detection of nuanced psychological threats is limited without advanced LLM reasoning.

---

## **6.2 Planned Improvements**

- Enable production-ready moderation API integration with rate limiting and caching
- Expand multilingual support (non-English violence pattern recognition)
- Add transformer-based contextual classifiers for deeper threat recognition
- Improve news detection with URL domain analysis and title heuristics
- Integrate active learning to continuously improve risk scoring
- Add visualization dashboards (risk trends, user graphs, post clusters)
- Migrate to async architecture for higher throughput