



Licence MIASHS deuxième année

Rapport de projet informatique

Détecteur de cyberharcèlement avancé en C assisté par IA

Projet réalisé du 1/01/2025 au 6/01/2025

Membres du groupe

DJEMATENE Dilan (44002835)
LACOMBE Ariane (42006077)

Dépôt GitHub :

<https://github.com/amazing-source/algo-preventioncyberharcelement>

Remerciements

Nous tenons à remercier nos enseignants pour cette proposition de projet stimulante, ainsi que la communauté open-source pour la documentation sur libcurl et json-c. Un remerciement particulier à Mistral AI pour l'accès gratuit à leur API, et à Claude (Anthropic) pour l'assistance au développement.

Table des matières

1	Introduction	5
1.1	Le cyberharcèlement : un fléau moderne aux conséquences dramatiques	5
1.2	Objectif du projet	5
1.3	Limites et considérations éthiques	5
2	Environnement de travail	5
2.1	Système d'exploitation	5
2.2	Outils de développement	6
2.3	Configuration de Code : :Blocks	6
3	Description du projet et objectifs	6
3.1	Fonctionnalités principales (v2.0)	6
3.2	Architecture à deux agents IA	6
4	Bibliothèques et technologies	6
4.1	Bibliothèque libcurl	6
4.2	Bibliothèque json-c	7
4.3	API Mistral AI	7
5	Utilisation de l'Intelligence Artificielle pour le développement	7
5.1	Méthodologie de développement assisté par IA	7
5.2	Bénéfices mesurables	7
6	Fonctionnement de l'IA intégrée (API Mistral)	7
6.1	Présentation de Mistral AI	7
6.2	Agent Générateur : Commentaires ultra-réalistes	7
6.3	Agent Analyseur : Analyse contextuelle avancée	8
6.4	Grille de scoring contextuelle	8
7	Travail réalisé	8
7.1	Fonctionnalités implémentées	8
7.1.1	Menu interactif principal	8
7.1.2	Mode 1 : Génération automatique	8
7.1.3	Mode 2 : Analyser un commentaire reçu	9
7.1.4	Mode 3 : Vérifier avant de publier	9
7.1.5	Mode 4 : Conversation	9
7.1.6	Mode 5 : Ressources d'aide	9
7.2	Modules du code source	9
8	Répartition du travail	9
9	Difficultés rencontrées	9
9.1	Gestion des accents sous Windows	9
9.2	Calibrage du scoring	10
9.3	Réalisme des commentaires générés	10

10 Bilan	10
10.1 Résultats obtenus	10
10.2 Conclusion	10
10.3 Perspectives	10
11 Bibliographie	11
12 Webographie	12
13 Annexes	13
A Cahier des charges	13
A.1 Objectifs fonctionnels	13
A.2 Contraintes techniques	13
B Exemple d'exécution du projet	13
C Manuel utilisateur	15
C.1 Installation	15
C.2 Configuration	15
C.3 Utilisation	15
D Structure du code source	15

1 Introduction

1.1 Le cyberharcèlement : un fléau moderne aux conséquences dramatiques

Le cyberharcèlement est devenu l'un des problèmes les plus graves de notre ère numérique. Contrairement au harcèlement traditionnel, il ne connaît ni frontières, ni horaires : la victime peut être atteinte 24 heures sur 24, 7 jours sur 7, dans l'intimité de son foyer.

Ces dernières années, nous avons assisté à l'émergence de phénomènes particulièrement alarmants. Des groupes organisés sur des plateformes comme Telegram orchestrent des campagnes de harcèlement de masse, ciblant des individus avec une violence inouïe. Ces "raids" coordonnés peuvent détruire la réputation d'une personne en quelques heures, diffuser des contenus intimes sans consentement, ou pousser des victimes dans leurs derniers retranchements.

Les conséquences du cyberharcèlement sont bien plus graves que ce que l'on imagine généralement. Au-delà de l'anxiété et de la dépression, le cyberharcèlement peut mener à l'irréparable. Chaque année, des victimes, souvent jeunes, mettent fin à leurs jours après avoir subi des campagnes de haine en ligne. Ces drames nous rappellent que derrière chaque écran se trouve un être humain, et que les mots peuvent tuer.

1.2 Objectif du projet

Ce projet est dédié à la lutte contre le cyberharcèlement. Notre objectif est de développer un outil capable de détecter automatiquement les commentaires toxiques, insultants ou menaçants, afin d'aider à la modération des espaces en ligne.

1.3 Limites et considérations éthiques

Il est important de souligner une réalité fondamentale : **il est extrêmement difficile de quantifier la gravité d'un propos par un simple nombre.** La perception d'un commentaire varie considérablement d'une personne à l'autre. Ce qui peut sembler anodin pour l'un peut être profondément blessant pour l'autre, en fonction de son vécu, de sa sensibilité, ou du contexte.

Notre système de scoring (0-100) ne prétend pas être une vérité absolue. Il s'agit d'une estimation basée sur des critères linguistiques et contextuels, inspirés des approches utilisées par les grandes plateformes (Google Perspective API, systèmes de modération de TikTok et Reddit). Nous avons fait de notre mieux pour créer un outil aussi précis et nuancé que possible, tout en reconnaissant ses limites inhérentes.

2 Environnement de travail

2.1 Système d'exploitation

Le développement s'est effectué sous Windows 11, avec MSYS2 MinGW-w64 pour disposer d'un environnement de compilation Unix-like compatible avec les bibliothèques nécessaires.

2.2 Outils de développement

- **IDE** : Code ::Blocks 20.03 configuré avec MinGW-w64
- **Compilateur** : GCC 13.2.0 (MSYS2 MinGW-w64)
- **Terminal** : MSYS2 MinGW64 pour l'installation des dépendances
- **Assistant IA** : Claude (Anthropic) pour le développement assisté
- **Gestionnaire de paquets** : pacman (MSYS2) pour libcurl et json-c

2.3 Configuration de Code ::Blocks

Configuration spécifique pour l'utilisation de libcurl et json-c :

- **Compiler search directories** : C :\msys64\mingw64\include
- **Linker search directories** : C :\msys64\mingw64\lib
- **Link libraries** : curl, json-c

3 Description du projet et objectifs

3.1 Fonctionnalités principales (v2.0)

1. **Menu interactif complet** avec 6 modes d'utilisation
2. **Génération automatique** de commentaires réalistes par IA
3. **Analyse de commentaires reçus** : pour vérifier si un message reçu est problématique
4. **Vérification avant publication** : pour tester son propre message avant de l'envoyer
5. **Mode conversation** : analyse multiple avec statistiques
6. **Ressources d'aide** : numéros d'urgence et sites de soutien (3018, 3114, etc.)
7. **Analyse contextuelle avancée** inspirée de Google Perspective API
8. **Conseils automatiques** : suggestions de reformulation si message problématique

3.2 Architecture à deux agents IA

Le système repose sur deux agents distincts :

- **Agent Générateur** : Produit des commentaires ultra-réalistes avec fautes d'orthographe, abréviations (mdr, tkt, jsp...), absence de majuscules, style authentique des réseaux sociaux
- **Agent Analyseur** : Évalue le niveau de toxicité selon des critères contextuels multiples

4 Bibliothèques et technologies

4.1 Bibliothèque libcurl

Utilisée pour les requêtes HTTP vers l'API Mistral. Permet la gestion des en-têtes d'authentification et l'envoi de requêtes POST contenant les prompts au format JSON.

4.2 Bibliothèque json-c

Essentielle pour le parsing des réponses de l'API. Permet d'extraire les champs score, catégorie et explication des réponses structurées de l'IA.

4.3 API Mistral AI

- **Authentification** : Clé API via header Authorization Bearer
- **Endpoint** : /v1/chat/completions
- **Modèle utilisé** : mistral-small-latest (gratuit)
- **Format** : Requêtes et réponses en JSON

5 Utilisation de l'Intelligence Artificielle pour le développement

5.1 Méthodologie de développement assisté par IA

Ce projet a été réalisé avec l'assistance de Claude (Anthropic), démontrant l'efficacité du pair programming humain-IA :

- **Debugging et résolution d'erreurs** : Identification rapide des problèmes de compilation et de logique
- **Architecture et design patterns** : Structure modulaire du code
- **Génération de code** : Modules de parsing JSON et communication API
- **Optimisation des prompts** : Affinement des instructions pour l'IA d'analyse
- **Documentation** : Aide à la rédaction du README et de ce rapport

5.2 Bénéfices mesurables

- **Réduction du temps de développement** : Environ 70%
- **Qualité du code** : Meilleures pratiques appliquées systématiquement
- **Apprentissage accéléré** : Compréhension approfondie de libcurl et des callbacks

6 Fonctionnement de l'IA intégrée (API Mistral)

6.1 Présentation de Mistral AI

Mistral AI est une entreprise française proposant des modèles de langage performants. L'API est accessible gratuitement avec le modèle `mistral-small-latest`.

6.2 Agent Générateur : Commentaires ultra-réalistes

Le prompt de génération a été conçu pour produire des commentaires authentiques :

```
1 Tu es un generateur de commentaires de reseaux sociaux
2 ULTRA REALISTES .
3 REGLES STRICTES pour etre realiste :
4 - Certains commentaires ont des fautes d'orthographe
5   (sa au lieu de ca, ses au lieu de c'est)
6 - Certains n'ont pas de majuscules du tout
7 - Certains utilisent des abbreviations
8   (tkt, mdr, ptdr, jsp, chui, pr, tt)
```

```

9 - Certains sont tres courts (1-5 mots)
10 genre 'nul' ou 'trop bien'
11 - Certains ont des emojis ou des ... ou ??? ou !!!
12 - Style VRAI Twitter/TikTok/YouTube francais

```

6.3 Agent Analyseur : Analyse contextuelle avancée

L'analyse ne se base jamais sur les mots isolés, mais sur le contexte global :

```

1 Tu es un detecteur de cyberharcelement base sur des
2 criteres contextuels (similaires a TikTok, Reddit
3 et Google Perspective API).
4
5 NE TE BASE JAMAIS sur les mots isoles.
6 Le mot "gros" peut etre amical ("ca va gros ?")
7 ou insultant ("sale gros").
8
9 Evaluate les elements suivants :
10 - Ton general : amical / neutre / moqueur / agressif
11 - Intention : humoristique, affective, insultante
12 - Structure : attaque directe ("t'es...")
13 - Cible : une personne specifique ou non
14 - Intensite : legere, moderee, forte
15 - Contexte linguistique : modificateurs aggravants

```

6.4 Grille de scoring contextuelle

- **0-20** : Sain / amical / neutre
- **20-40** : Rude / désagréable mais non insultant
- **40-60** : Insulte légère / moquerie isolée
- **60-80** : Attaque personnelle claire / humiliation
- **80-100** : Haine, body-shaming agressif, menace, discrimination

7 Travail réalisé

7.1 Fonctionnalités implémentées

7.1.1 Menu interactif principal

```

1 -----
2           MENU PRINCIPAL
3 -----
4 [1] Generation automatique (IA)
5 [2] Analyser un commentaire recu
6 [3] Verifier avant de publier
7 [4] Mode conversation (multi-analyse)
8 [5] Ressources d'aide
9 [6] Quitter
10 -----

```

7.1.2 Mode 1 : Génération automatique

L'utilisateur choisit un thème et un nombre de commentaires. L'IA génère des commentaires ultra-réalistes (avec fautes, abréviations, style internet) puis les analyse un

par un avec statistiques finales.

7.1.3 Mode 2 : Analyser un commentaire reçu

Destiné aux **victimes potentielles**. L'utilisateur colle un message qu'il a reçu pour vérifier s'il est problématique. Si le score est élevé, le programme affiche automatiquement les ressources d'aide (3018, netecoute.fr).

7.1.4 Mode 3 : Vérifier avant de publier

Destiné aux **auteurs**. Permet de tester son propre message avant de l'envoyer. Le programme affiche des avertissements graduels selon la gravité et propose des conseils de reformulation si le message est problématique.

7.1.5 Mode 4 : Conversation

Permet d'analyser plusieurs commentaires à la suite. Affiche des statistiques globales à la fin (score moyen, répartition par catégorie, commentaire le plus toxique).

7.1.6 Mode 5 : Ressources d'aide

Affiche les numéros et sites d'aide :

- 3018 : Numéro national contre le cyberharcèlement
- 3114 : Prévention du suicide (24h/24)
- netecoute.fr, e-enfance.org, internet-signalement.gouv.fr

7.2 Modules du code source

- **main.c** : Menu interactif et logique principale
- **api_client.c/.h** : Communication HTTP avec l'API Mistral
- **comment_generator.c/.h** : Agent 1 - Génération réaliste
- **harassment_detector.c/.h** : Agent 2 - Analyse contextuelle
- **stats.c/.h** : Calcul et affichage des statistiques
- **config.h** : Configuration de la clé API

8 Répartition du travail

- **DJEMATENE Dilan** : Architecture générale, développement du code C, intégration de l'API Mistral.
- **LACOMBE Ariane** : Recherches sur le cyberharcèlement, rédaction du rapport, tests utilisateur, documentation, relecture, tests et debugging, conception des prompts IA

9 Difficultés rencontrées

9.1 Gestion des accents sous Windows

L'affichage des caractères accentués dans la console Windows a nécessité une configuration spéciale avec `SetConsoleOutputCP(65001)` pour passer en UTF-8.

9.2 Calibrage du scoring

Trouver le bon équilibre entre faux positifs et faux négatifs a demandé de nombreuses itérations sur les prompts. L'approche contextuelle (ne pas se baser sur les mots isolés) a grandement amélioré la précision.

9.3 Réalisme des commentaires générés

Les premiers commentaires générés étaient trop "propres" et ne ressemblaient pas à de vrais commentaires internet. L'ajout d'instructions spécifiques (fautes, abréviations, style oral) a résolu ce problème.

10 Bilan

10.1 Résultats obtenus

Le détecteur v2.0 offre une expérience interactive complète avec trois modes d'utilisation. L'analyse contextuelle permet une évaluation nuancée qui distingue les usages amicaux des usages insultants d'un même mot.

10.2 Conclusion

Ce projet démontre qu'il est possible de créer un outil de détection de cyberharcèlement en C, exploitant l'intelligence artificielle pour une analyse sémantique fine. Bien que le système ne soit pas parfait (la perception de la gravité reste subjective), il constitue une base solide pour aider à la modération des espaces en ligne.

La lutte contre le cyberharcèlement est l'affaire de tous. Nous espérons que ce projet contribuera, à son échelle, à sensibiliser et à protéger les victimes potentielles.

10.3 Perspectives

- Intégration avec des APIs de réseaux sociaux (Discord, Reddit)
- Interface graphique pour une utilisation grand public
- Système d'apprentissage pour améliorer la précision au fil du temps
- Extension multilingue

11 Bibliographie

[1] Aucun ouvrage papier consulté.

12 Webographie

- [CURL] Documentation libcurl : <https://curl.se/libcurl/c/>
- [JSONC] Documentation JSON-C : <https://github.com/json-c/json-c>
- [MISTRAL] Documentation API Mistral : <https://docs.mistral.ai/>
- [PERSPECTIVE] Google Perspective API : <https://perspectiveapi.com/>
- [MSYS2] MSYS2 guide d'installation : <https://www.msys2.org/>
- [CLAUDE] Assistant IA Claude : <https://claude.ai/>

13 Annexes

Annexe A : Cahier des charges

Annexe A.1 : Objectifs fonctionnels

- Menu interactif avec plusieurs modes
- Générer des commentaires réalistes via l'Agent 1
- Permettre l'analyse de commentaires personnalisés
- Analyser chaque commentaire via l'Agent 2 (analyse contextuelle)
- Attribuer un score de toxicité (0-100)
- Classifier selon les catégories définies
- Fournir une explication pour chaque classification
- Générer des statistiques globales

Annexe A.2 : Contraintes techniques

- Langage C uniquement
- Compilation avec GCC sans warnings
- Exécution en ligne de commande avec menu interactif
- Communication API via libcurl
- Parsing JSON via json-c

Annexe B : Exemple d'exécution du projet

```
1 $ ./cyberharcelement_detector
2
3 =====
4     DETECTEUR DE CYBERHARCELEMENT v2.0
5         Analyse & Prevention par Intelligence Artificielle
6 =====
7
8 -----
9         MENU PRINCIPAL
10 -----
11    [1] Generation automatique (IA)
12    [2] Analyser un commentaire recu
13    [3] Verifier avant de publier
14    [4] Mode conversation (multi-analyse)
15    [5] Ressources d'aide
16    [6] Quitter
17 -----
18 Votre choix : 5
19
20 =====
21         RESSOURCES D'AIDE ET DE PREVENTION
22 =====
23
24 Si vous etes VICTIME de cyberharcelement :
25 -----
26     - 3018 : Numéro national contre le cyberharcelement (gratuit)
27     - https://www.netecoute.fr : Ecoute et conseils
```

```
29 - https://www.e-enfance.org : Protection des mineurs  
30  
31 Si vous avez des PENSEES SOMBRES :  
-----  
33 - 3114 : Numéro national de prévention du suicide (24h/24)  
34 - https://www.sos-amitie.com : Ecoute 24h/24
```

```
1 Votre choix : 2  
2  
3 === ANALYSER UN COMMENTAIRE REÇU ===  
4 (Pour vérifier si un message que vous avez reçu est problématique)  
5  
6 Collez le commentaire à analyser :  
7 > ca va gros ? t'as passé un bon weekend ?  
8  
9 [Analyse en cours...]  
10  
11 --- RESULTAT DE L'ANALYSE ---  
12 | Score : 8/100  
13 | Catégorie : sain  
14 |-----  
15 | Explication :  
16 | Le terme "gros" est utilisée de manière amicale  
17 | et familière. Le ton est bienveillant.  
18 |-----  
19 [OK] Ce commentaire semble sain.
```

```
1 Votre choix : 3  
2  
3 === VERIFIER AVANT DE PUBLIER ===  
4 (Testez votre message AVANT de l'envoyer)  
5  
6 Ecrivez votre commentaire :  
7 > t'es vraiment qu'un gros débile  
8  
9 [Vérification en cours...]  
10  
11 #####  
12 # VÉRIFICATION AVANT PUBLICATION #  
13 #####  
14  
15 [XXX] STOP ! Ce commentaire est une ATTAQUE.  
16 Catégorie : humiliation (score: 72/100)  
17  
18 -> Ceci peut constituer du HARCELEMENT.  
19 -> La victime pourrait PORTER PLAINE.  
20 -> Vous risquez des POURSUITES JUDICIAIRES.  
21  
22 NE PUBLIEZ PAS CE MESSAGE.  
23 Respirez. Fermez l'écran. Revenez plus tard.  
24  
25 #####  
26  
27 Voulez-vous des suggestions pour reformuler ? (o/n) : o  
28  
29 Conseils pour reformuler :  
30 - Exprimez votre désaccord sans attaquer la personne  
31 - Utilisez "je pense que..." plutôt que "tu es..."  
32 - Critiquez les idées, pas les individus
```

Annexe C : Manuel utilisateur

Annexe C.1 : Installation

```
1 # Ouvrir MSYS2 MinGW 64-bit
2
3 # Installation des dependances
4 pacman -S mingw-w64-x86_64-curl
5 pacman -S mingw-w64-x86_64-json-c
6
7 # Compilation
8 make
```

Annexe C.2 : Configuration

Modifier le fichier config.h avec votre clé API Mistral (la mienne est mise à disposition) :

```
1 #define API_KEY "votre-cle-mistral-ici"
```

Annexe C.3 : Utilisation

```
1 ./cyberharcelement_detector
```

Puis suivre le menu interactif.

Annexe D : Structure du code source

```
1 projet/
2   |-- main.c           # Menu interactif
3   |-- api_client.c     # Communication API
4   |-- api_client.h
5   |-- comment_generator.c # Agent 1
6   |-- comment_generator.h
7   |-- harassment_detector.c # Agent 2
8   |-- harassment_detector.h
9   |-- stats.c           # Statistiques
10  |-- stats.h
11  |-- config.h          # Configuration
12  |-- Makefile
13  |-- README.md
14  |-- rapport/
15    |-- main.tex
16    |-- pagedegarde.sty
17    |-- bordure.png
18    |-- logo_Paris_Nanterre_couleur_RVB.png
```