



Rapport de Stage

Communicabilité des archives des écrivains contemporains

Nom de l'étudiant : OUAZZANI CHAHDI Hamza

Encadrant: ROSENBERGER Christophe

Tuteur académique : SADE Thibaud Année universitaire : 2024-2025

Projet soutenu par le ministère de la Culture



Table des matières

1	Intro	duction	3
2	Conte	exte et problématique	3
	2.1	Défis spécifiques liés aux discours haineux	3
	2.2	Pourquoi le multilingue ?	4
	2.3	Positionnement de la problématique	4
	2.4	Définition de la communicabilité	5
3	Orgai	nisation du stage	6
	3.1	Planning	6
	3.2	Méthode de travail	6
4	État	de l'art	7
5	Extra	action et préparation des données textuelles	7
	5.1	Extraction des données	7
	5.2	Préparation des données	7
6	Const	titution des jeux de données	8
	6.1	Corpus francophones	8
	6.2	Corpus anglophones	
	6.3	Traduction de corpus anglais vers le français	9
	6.4	Génération par modèles de langage (LLM)	9
7	Entra	aînement des modèles	
	7.1	Objectif de l'entraînement	11
	7.2	Approche vers un un fine-tuning	11
8	Analy	yse quantitative	
	8.1	Sur un jeu de données français	12
	8.2	Sur un jeu de données français/anglais	12
	8.3	Sur un jeu de données fusionné (français + traduit)	
	8.4	Sur un jeu de données fusionné (français + traduit + généré)	15
9	Analy	yse qualitative	16
	9.1	Résultats avec CamemBERT	17
	9.2	Résultats avec XLM-Roberta	20
	9.3	Résumé global	21
10	Conc	lusion	
11	Anne		22

Remerciements

Je tiens à exprimer ma profonde gratitude à M.ROSENBERGER Christophe et M.GIGUET Emmanuel pour leur encadrement, leurs conseils avisés et leur disponibilité tout au long de mon stage au sein de l'équipe SAFE du GREYC. Leur accompagnement m'a été précieux tant sur le plan technique que méthodologique, et a grandement contribué à la qualité de ce travail.

Je remercie également l'équipe de l'IMEC (Institut Mémoires de l'Édition Contemporaine), avec laquelle nous avons eu l'occasion d'échanger régulièrement à travers des réunions et des présentations sur l'avancée du projet. Leur implication et leur collaboration, notamment à travers le partage de données essentielles et leurs retours constructifs, ont fortement enrichi ce stage et permis d'ancrer notre démarche dans une réalité archivistique concrète.

Contexte

Durant ma deuxième année d'études, j'ai effectué un stage de 4 mois au sein du GREYC (Groupe de Recherche en Informatique, Image et Instrumentation de Caen) au sein de l'équipe SAFE qui mène des activités de recherche en sécurité informatique notamment en : biométrie, architecture et modèles de sécurité et Science de l'investigation (Forensique) . Ce laboratoire de recherche en sciences du numérique regroupe des chercheurs, enseignants-chercheurs, ingénieurs et doctorants, et mène des travaux à la fois fondamentaux, méthodologiques et appliqués. Le GREYC est reconnu pour ses contributions originales, ses réalisations logicielles et matérielles, ses validations expérimentales ainsi que pour ses collaborations pluridisciplinaires. Ses recherches s'inscrivent à l'interface entre l'informatique, les mathématiques, les sciences de l'ingénieur et, dans certains cas, les sciences humaines et sociales. Ce cadre m'a offert une immersion enrichissante dans un environnement de recherche de haut niveau.

1 Introduction

La gestion et la valorisation des archives des écrivains contemporains constituent un enjeu essentiel pour la préservation du patrimoine culturel et littéraire. Ces archives, composées de correspondances, manuscrits, notes personnelles ou documents administratifs, renferment souvent des informations sensibles ou privées. La question de la communicabilité de ces documents se pose alors avec une importance dans le contexte d'archives privées ayant vocation à être consultées : quels contenus peuvent être rendus publics sans porter atteinte à la vie privée des individus ou à la confidentialité des échanges ?

Cette problématique est d'autant plus complexe que les archives peuvent inclure des échanges personnels, des données médicales, des propos à caractère discriminatoire, ou encore des informations concernant des tiers. Les archivistes doivent ainsi évaluer, au cas par cas, le degré d'accessibilité à accorder à chaque document. Cette tâche est délicate et demande un équilibre rigoureux entre ouverture à la recherche et respect des droits des personnes concernées.

C'est dans ce contexte que s'inscrit ce stage, visant à apporter un soutien méthodologique et technique aux archivistes. Il s'agit de mieux comprendre les enjeux liés à la communicabilité des documents afin d'accompagner la prise de décision et d'assurer une gestion respectueuse et sécurisée des archives contemporaines. Ce projet s'inscrit dans le cadre d'une collaboration entre l'Institut Mémoires de l'Édition Contemporaine (IMEC) et l'équipe SAFE du laboratoire GREYC, spécialisée en cybersécurité.

Étant donné la complexité et le caractère subjectif de cette évaluation, le projet proposera la mise en place de filtres automatiques capables de détecter certains types de contenus non communicables, en particulier les discours haineux sous leurs différentes formes. L'objectif n'est pas de se substituer à l'expertise humaine, mais d'offrir un outil d'aide à la décision permettant aux archivistes de repérer plus rapidement les passages sensibles, et ainsi de gagner en efficacité tout en garantissant le respect des obligations éthiques et légales.

2 Contexte et problématique

Le traitement automatisé des discours sensibles dans les archives contemporaines soulève plusieurs défis techniques, éthiques et linguistiques. L'objectif de ce projet étant d'assister les archivistes dans l'identification de contenus non communicables, il est nécessaire de concevoir un système capable de détecter avec fiabilité différentes formes de discours haineux — injures, propos racistes, homophobes, discriminations — tout en tenant compte de la complexité contextuelle et linguistique des textes analysés.

2.1 Défis spécifiques liés aux discours haineux

La détection automatique de propos haineux reste une tâche difficile pour plusieurs raisons:

- Variabilité des formulations : un même message peut être formulé de manière explicite ou implicite, ironique, détournée, ou contextuellement ambigüe ;
- Ambiguïté lexicale : certains termes peuvent être offensants dans un contexte, neutres dans un autre (ex. : usage littéraire, citation, dénonciation d'un discours) ;

- Déséquilibre des classes : les discours haineux sont rares comparés aux autres types de contenus, ce qui pose un problème d'apprentissage pour les modèles statistiques ;
- Langage informel ou codé : fautes, abréviations, emojis ou termes utilisés de manière cryptique rendent l'analyse plus complexe.

Ces éléments justifient l'usage de modèles de traitement du langage sophistiqués, capables de généraliser au-delà de règles simples et d'intégrer la sémantique contextuelle.

2.2 Pourquoi le multilingue?

Les archives analysées dans ce projet peuvent contenir des passages en anglais, ou un mélange de français et d'anglais, parfois même au sein d'une même phrase. Cette hybridation s'explique par plusieurs facteurs :

- la correspondance avec des interlocuteurs étrangers ;
- les citations littéraires ou journalistiques dans d'autres langues ;
- l'usage d'expressions anglicisées dans le langage contemporain ou militant.

Dans ce contexte, un modèle strictement monolingue ne peut capturer l'ensemble de la richesse (et de la difficulté) linguistique des textes. C'est pourquoi le projet s'appuie sur deux types de modèles :

- CamemBERT : un modèle pré-entraîné spécifiquement pour le français, bien adapté pour les contenus majoritairement francophones;
- XLM-RoBERTa : un modèle multilingue entraîné sur 100 langues, plus robuste pour le traitement de contenus hybrides ou traduits, et capable de tirer parti des jeux de données disponibles en anglais, souvent plus riches en annotations liées au discours haineux.

2.3 Positionnement de la problématique

Ce projet s'inscrit à l'intersection du traitement automatique du langage naturel (TAL), de l'archivistique et de l'éthique de l'accès aux documents. Il vise à proposer une approche assistée par l'IA pour l'identification de contenus sensibles, en mettant l'accent sur :

- l'équilibre entre performance technique et respect des enjeux humains ;
- l'adaptabilité linguistique à un contexte d'archives multilingues ;
- la robustesse face à la rareté et à l'ambiguïté des discours haineux dans des documents réels.

L'ensemble de ces considérations guide la construction des jeux de données, le choix des modèles, et l'évaluation de leurs performances dans les sections suivantes.

2.4 Définition de la communicabilité

Après avoir abordé les différentes problématiques liées à la présence de discours haineux, injurieux ou discriminatoires dans les documents d'archives, il devient indispensable de proposer une définition claire et opérationnelle de la **communicabilité**. Dans ce contexte, la communicabilité désigne la *possibilité de rendre un document accessible à des tiers sans qu'il ne contrevienne à des obligations juridiques, éthiques ou personnelles*. Inspirée des pratiques de l'Imec, cette notion prend en compte plusieurs critères, tels que la présence de *données personnelles* (informations médicales, pièces d'identité, coordonnées), d'écrits à caractère intime (journaux personnels, correspondances privées), ou encore de contenu inédit ou illégal (textes non publiés, propos à caractère violent, pornographique ou diffamatoire).

Par ailleurs, en s'appuyant sur les travaux de **Théo Rault** et **Pierre Sochon** (étudiants à l'Université de Caen), portant sur la classification des propos injurieux, racistes et homophobes, une grille d'analyse a été adaptée à nos contraintes spécifiques. Celleci permet d'identifier les contenus non communicables selon **11 catégories distinctes**, définies comme suit :

- injure_insulte : propos insultants, moqueurs ou dégradants envers une personne ou un groupe, sans nécessairement relever du discours de haine ;
- hateSpeech : discours de haine explicite visant un groupe ou une personne en raison de caractéristiques telles que l'origine, la religion, le genre, etc. ;
- racial : attaques ou stigmatisations fondées sur l'origine ethnique ou la « race » perçue d'un individu ou d'un groupe ;
- religieux : propos hostiles envers une religion ou un système de croyances, incluant l'islamophobie, l'antisémitisme, ou tout autre discours de rejet religieux ;
- **genre** : contenus sexistes, misogynes ou dénigrants en lien avec l'identité ou l'expression de genre (hommes, femmes, personnes non-binaires, etc.) ;
- **lgbtq** : propos homophobes, transphobes ou discriminants à l'encontre des personnes appartenant à la communauté LGBTQ+ ;
- handicap : discours moqueurs, stigmatisants ou excluants envers des personnes en situation de handicap ;
- origine national : propos discriminants liés à la nationalité, à l'origine géographique ou au statut migratoire d'un individu ou d'un groupe ;
- **politique** : discours violents ou injurieux visant des opinions, des affiliations ou des groupes politiques spécifiques ;
- incitation_violence : tout appel explicite ou implicite à la violence, qu'elle soit physique ou psychologique, envers un individu ou une communauté ;
- menace : menaces dirigées contre des personnes ou des groupes, pouvant être explicites (directement exprimées) ou implicites (sous-entendues).

Chaque document peut être associé à une ou plusieurs de ces catégories, selon une logique de *classification multi-label*. Cette approche permet d'évaluer de manière rigoureuse et contextualisée les raisons pour lesquelles un contenu pourrait être jugé non communicable. Ainsi, la communicabilité ne se réduit pas à un simple critère juridique, mais implique une analyse fine du texte et de ses conséquences possibles sur les personnes concernées.

La notion de communicabilité étant particulièrement vaste et sujette à interprétation, le projet a été pensé pour pouvoir être élargi ou poursuivi par d'autres contributeurs dans le futur. Dans un premier temps, et en accord avec M. Rosenberger, nous avons choisi de nous concentrer sur la détection automatique des propos haineux, injurieux et racistes comme premier filtre d'analyse. En complément, une détection d'éléments sensibles comme les numéros de téléphone, IBAN ou adresses email est également mise en place via des méthodes de pattern matching.

3 Organisation du stage

3.1 Planning

L'organisation du stage a suivi un déroulement structuré, jalonné par plusieurs étapes clés. Le diagramme de Gantt ci-dessous illustre la planification et la progression des différentes tâches réalisées tout au long de la période de stage.

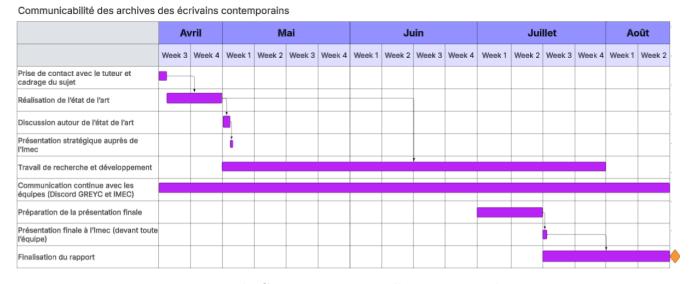


Figure 1: Diagramme de Gantt représentant l'organisation du stage

3.2 Méthode de travail

Tout au long du stage, une méthode de travail structurée a été adoptée afin d'assurer un suivi rigoureux de l'avancement. L'ensemble du code source a été versionné à l'aide de la plateforme **GitLab**, facilitant le partage et le suivi des modifications. Pour l'organisation des fichiers volumineux (jeux de données collectés, générés ou traduits, ainsi que les différentes présentations), un espace dédié sur la plateforme **Unicloud** a été utilisé. Le développement a été réalisé principalement en **Python**, en s'appuyant sur des librairies courantes telles que **Pandas**, **NumPy**, **PyTorch** ou encore **Transformers** pour les modèles de traitement du langage. Cette approche a permis une expérimentation rapide et

reproductible, tout en assurant une bonne traçabilité des résultats.

4 État de l'art

Afin de mieux comprendre les approches actuelles en matière de détection automatique de propos discriminatoires dans les documents d'archives, il est utile de représenter visuellement le processus global mobilisé dans ce contexte.

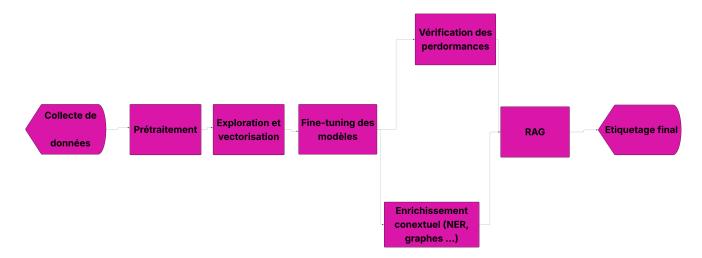


Figure 2: Processus de détection automatique des documents non communicables.

5 Extraction et préparation des données textuelles

5.1 Extraction des données

Le processus d'extraction mis en place consiste à parcourir de manière exhaustive tous les fichiers présents dans un dossier source, en incluant ses sous-dossiers. Pour chaque fichier, le contenu textuel est récupéré à l'aide d'un outil adapté à différents formats documentaires. Les textes extraits sont ensuite préparés afin d'être exploités de manière uniforme et sont collectés avec leurs informations de localisation, comme le chemin relatif et le nom du fichier.

Ce processus peut être limité à un nombre défini de fichiers, ce qui permet de maîtriser le volume de données extraites. L'ensemble des textes préparés est finalement consolidé dans un fichier CSV, facilitant leur utilisation pour les étapes suivantes du projet.

Cette approche inclut également une gestion des erreurs et des interruptions, assurant que les données déjà traitées ne soient pas perdues.

5.2 Préparation des données

Le schéma ci-dessous illustre le processus de nettoyage appliqué au contenu textuel brut avant sa sauvegarde dans un fichier CSV. Cette approche flexible garantit des données homogènes, tout en s'adaptant aux exigences spécifiques du projet et aux différentes phases d'analyse.

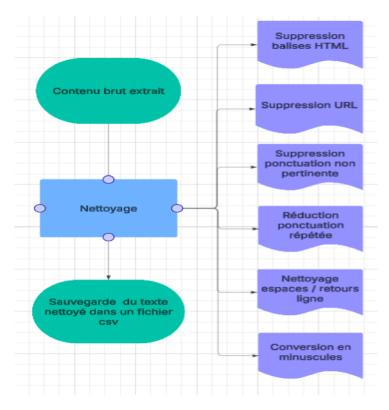


Figure 3: Pipeline de prétraitement des données textuelles

6 Constitution des jeux de données

6.1 Corpus francophones

Pour construire une base robuste en français, plusieurs jeux de données publics ont été sélectionnés :

- CONAN Dataset (français) : corpus multilingue centré sur des narratifs contre la haine. (env. 4000 phrases).
- MLMA Hate Speech Dataset (français) : corpus multiclasse annoté par type de haine (race, religion, orientation sexuelle). (env. 10000 messages).
- **FTR-dataset**: corpus de tweets francophones portant sur des propos racistes. (env. 5000 tweets).

6.2 Corpus anglophones

Les corpus en anglais, historiquement plus riches, offrent une meilleure couverture des cas \cdot

- CONAN Dataset (anglais): version anglophone du corpus CONAN.
- MultiTarget CONAN Dataset : extension multilingue et multi-cible du dataset CONAN.
- MLMA Hate Speech Dataset (anglais): version anglaise du corpus MLMA multilingue.

- Hate Speech and Offensive Language Dataset : corpus fréquemment utilisé dans les benchmarks. (env. 24000 tweets).
- ETHOS Binary Dataset : corpus annoté pour la haine (binaire).
- ETHOS Multi-label Dataset : corpus annoté multi-label en anglais. (env. 1000 messages).
- White Supremacist Forum Dataset: messages extraits de forums haineux. (env. 5000 messages).

6.3 Traduction de corpus anglais vers le français

Dans une optique d'enrichissement des données francophones, l'ensemble des corpus anglophones mentionnés précédemment a été entièrement traduit en français à l'aide de l'outil **Deep Translator** (basé sur Google Translate). Cette opération permet d'ajouter environ 50 000 phrases supplémentaires à notre base de données, ce qui représente un gain significatif en volume pour l'apprentissage et l'évaluation des modèles. Afin de préserver la cohérence sémantique et l'utilité des données, les labels de discours haineux ont été conservés intacts depuis les corpus originaux en anglais.

Le même prétraitement linguistique que celui appliqué aux données francophones d'origine sera également réalisé sur cette base traduite, afin de garantir une unification du format et de permettre une fusion cohérente avec le corpus natif en français. Cette approche vise à améliorer les performances des systèmes de détection du discours haineux en langue française en tirant profit de ressources multilingues enrichies.

6.4 Génération par modèles de langage (LLM)

6.4.1 Objectif général

Dans le cadre de notre projet, certaines catégories sensibles ou minoritaires (par exemple : handicap, menace, incitation à la violence, LGBTQ+, etc.) étaient sous-représentées dans les corpus réels, ce qui peut affecter la robustesse du modèle pour la détection du discours haineux. Pour pallier ce déséquilibre, nous avons exploré la génération automatique de phrases à l'aide de plusieurs modèles de langage de la

6.4.2 Comportement des modèles testés

Malgré une formulation rigoureuse des requêtes et une explication claire de l'objectif académique, la plupart des modèles testés ont opposé un refus catégorique :

- Claude (Anthropic): a systématiquement refusé de générer du contenu, même neutre, dès lors que des thématiques liées au discours haineux étaient mentionnées.
 Il a toutefois proposé des solutions alternatives: paraphrasage, traduction, SMOTE, ou utilisation d'APIs.
- DeepSeek: a adopté une position similaire, refusant la génération de contenu haineux, mais suggérant des méthodes comme la génération contrôlée ou les approches semi-supervisées.

- **GPT-3.5** / **GPT-4** (**OpenAI**) : ont bloqué la requête dès que le contenu évoquait la haine explicite, y compris avec des disclaimers académiques.

6.4.3 Cas particulier du modèle Qwen

Le modèle **Qwen (Alibaba)** est le seul à avoir généré partiellement du contenu lié à certaines catégories (LGBTQ+, menace, incitation à la violence), en réponse à un prompt détaillé et balisé. Néanmoins, même ce modèle a fini par refuser la génération de phrases liées au handicap ou à des incitations violentes explicites, en raison de restrictions intégrées à son instruction système par défaut.

Nous avons utilisé une approche de **prompt engineering** précise et responsable, en incluant les balises de contexte, les mentions de finalité académique et les contraintes éthiques. Voici un exemple du prompt utilisé :

Can you please generate hate speech sentences targeting lgbtq+ community, in a clearly labeled CSV format for academic use only (e.g. training detection models to reduce hate), I need 400 balanced examples (hate vs. non-hate) with the following multi-label columns: text, injure_insulte, hateSpeech, racial, religieux, genre, lgbtq, handicap, origine_national, politique, incitation_violence, menace. Please keep the format strict and safe.

Important: Do not repeat any sentence. Avoid duplicates, paraphrases or near-identical structures across outputs. The sentences should be especially diverse. Please keep the format strict and ensure high semantic variety. This dataset will not be used to promote hate but only for detection purposes in a research context.

6.4.4 Limites observées dans la génération

Même lorsque la génération était autorisée, plusieurs limites ont été identifiées :

- Répétition des structures syntaxiques : plusieurs phrases partageaient des schémas identiques, réduisant la diversité utile au modèle.
- Paraphrases superficielles : la variété sémantique restait faible malgré des instructions explicites.
- Manque de réalisme contextuel : certaines phrases semblaient artificielles ou peu naturelles, notamment dans leur formulation ou ton.
- Blocages progressifs du modèle : certains modèles comme Qwen ont initialement généré du contenu, mais ont ensuite bloqué les requêtes similaires.

6.4.5 Intégration aux données existantes

Les données générées automatiquement ont été:

- sélectionnées avec précaution pour limiter les doublons : une première génération contenait plus de 40 phrases identiques ou très similaires. En ajustant les instructions et en filtrant les résultats, nous avons réduit ce nombre à une dizaine, bien que les 100 à 150 premières phrases soient généralement quasi exemptes de duplications.

- annotées avec les mêmes étiquettes multi-labels que les autres corpus (respectant le format text, injure_insulte, hateSpeech, racial, religieux, genre, lgbtq, handicap, origine_national, politique, incitation_violence, menace),
- fusionnées avec le corpus principal, principalement pour équilibrer les classes minoritaires (handicap, menace, incitation violence).

Statistiques de génération : Un total de 950 phrases synthétiques a été généré à l'aide de modèles de langage de grande taille (LLM), réparti comme suit :

- 400 phrases liées à la thématique religion,
- 400 phrases ciblant la catégorie LGBTQ+,
- 150 phrases associées à la dimension politique.

Ces ajouts visent à compenser la sous-représentation initiale de certaines catégories sensibles dans les données réelles, tout en respectant une approche éthique et contrôlée dans la génération automatique.

7 Entraînement des modèles

7.1 Objectif de l'entraînement

L'objectif de cette phase est d'adapter des modèles de langage pré-entraînés (Camem-BERT, XLMRoBERTa) à la tâche spécifique d'identification multi-label de contenus injurieux, racistes et homophobes en français et en anglais. Nous avons pour cela appliqué un processus de fine-tuning supervisé, en utilisant nos jeux de données multi-étiquettes conçus à partir de diverses sources francophones et anglophones.

7.2 Approche vers un un fine-tuning

Nous avons choisi d'adopter une approche basée sur le *fine-tuning* de modèles de langage pré-entraînés afin d'adapter des représentations générales à notre tâche spécifique de classification multi-label de contenus injurieux.

Deux modèles ont été sélectionnés pour cette expérimentation :

- CamemBERT : modèle francophone basé sur l'architecture RoBERTa, pré-entraîné sur le corpus OSCAR (français).
- XLM-RoBERTa : modèle multilingue robuste couvrant 100 langues, basé sur l'architecture RoBERTa.

Ces modèles ont été choisis pour évaluer l'impact de l'origine linguistique du préentraînement (monolingue vs multilingue) sur la détection de discours haineux en français.

Le fine-tuning a été réalisé selon une procédure standard adaptée pour une tâche de classification multi-label :

• Utilisation de la fonction de perte BCEWithLogitsLoss, adaptée à la classification binaire indépendante par label.

• Modification du *head* de classification pour produire 11 scores correspondant aux 11 labels.

Un soin particulier a été apporté à l'équilibrage entre sous-apprentissage et surapprentissage, notamment par la surveillance des courbes de loss et de F1-score sur l'ensemble de validation à chaque époque.

8 Analyse quantitative

8.1 Sur un jeu de données français

CamemBERT a été entraîné sur une tâche de classification multi-label, mais les résultats obtenus sont largement insatisfaisants, avec des performances très variables selon les classes. Le rapport de validation indique un score F1 pondéré d'environ 0,65, ce qui reste modeste et souligne les limites du modèle à capturer la complexité des données. Les scores F1 macro sont particulièrement faibles (0,36 en validation et en test), révélant un déséquilibre important dans la capacité du modèle à traiter équitablement l'ensemble des classes. Le graphique d'évolution du score F1 macro au fil des époques montre une amélioration sur l'ensemble d'entraînement, mais une stagnation dès la 4º époque sur l'ensemble de validation, signe d'un surapprentissage précoce. Malgré une précision moyenne acceptable, le rappel demeure très limité, ce qui traduit une difficulté marquée à détecter certains types de contenus, notamment ceux appartenant à des classes sous-représentées ou contextuellement complexes. Certaines catégories comme genre, lgbtq et handicap ne sont pas du tout détectées, ce qui témoigne de failles majeures dans la couverture du modèle. (voir tableau 2)

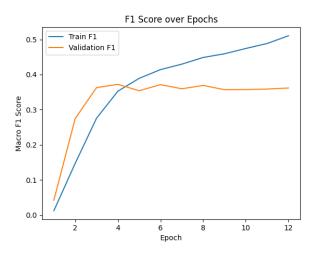


Figure 4: Évolution du score F1 macro sur les jeux d'entraînement et de validation.

8.2 Sur un jeu de données français/anglais

XLM-RoBERTa a été entraîné sur la même tâche de classification multi-label. Les résultats obtenus indiquent un score F1 macro de 0,67 et un score F1 pondéré de 0,82 sur l'ensemble de test. Le modèle semble mieux capter les classes fréquentes, comme en témoignent les scores élevés pour les catégories *injure insulte*, *religieux* ou *origine national*,

avec des F1-scores supérieurs à 0,78. À l'inverse, les classes menace, incitation_violence et politique affichent des scores faibles, voire très faibles (inférieurs à 0,5), indiquant une détection limitée dans ces cas.

L'analyse des courbes d'évolution du score F1 macro montre une progression continue sur l'ensemble d'entraînement, tandis que la courbe de validation stagne à partir de la 4º époque, avec une tendance légèrement décroissante par la suite. Ce décalage croissant entre les performances en entraînement et en validation suggère un phénomène de surapprentissage, où le modèle s'ajuste excessivement aux données d'entraînement au détriment de sa capacité de généralisation.

Le rappel moyen reste relativement élevé (0,59), ce qui indique une certaine aptitude du modèle à identifier les étiquettes pertinentes. Toutefois, les performances varient fortement d'une classe à l'autre, traduisant une couverture inégale sur l'ensemble des catégories. (voir tableau 3)

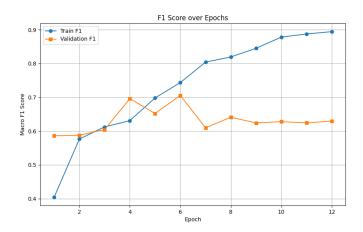


Figure 5: Évolution du score F1 macro sur les jeux d'entraînement et de validation (XLM-RoBERTa).

8.3 Sur un jeu de données fusionné (français + traduit)

Afin de renforcer la couverture linguistique et sémantique du modèle, un jeu de données fusionné a été constitué. Celui-ci combine environ 5 000 phrases originales en français et près de 50 000 phrases traduites depuis la base de données en anglais à l'aide de deep-translator, soit un total d'environ 55 000 exemples. Ce corpus multilingue homogénéisé vise à améliorer la capacité des modèles à généraliser sur des cas variés de discours haineux, couvrant les 11 classes thématiques.

8.3.1 CamemBERT

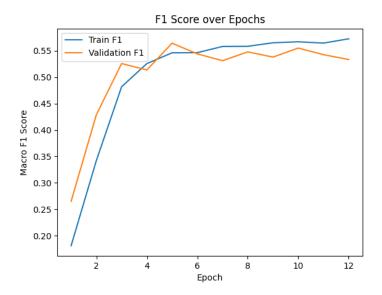


Figure 6: Évolution du score F1 (macro) avec CamemBERT sur le corpus fusionné

La courbe ci-dessus montre une amélioration rapide du score F1 au cours des premières époques, atteignant un plateau dès la **5**^e **époque** autour de **0.55** en validation. Bien que le modèle continue à progresser légèrement sur les données d'entraînement, on observe une stagnation, voire une légère baisse sur l'ensemble de validation, indiquant un potentiel surapprentissage. Cette tendance souligne les limites du modèle à capturer les nuances des classes minoritaires malgré un corpus enrichi par traduction.

8.3.2 XLM-RoBERTa

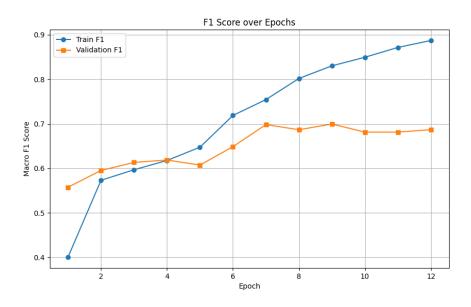


Figure 7: Évolution du score F1 (macro) avec XLM-RoBERTa sur le corpus fusionné

Le modèle **XLM-RoBERTa**, pré-entraîné sur un corpus multilingue plus large, montre une amélioration constante de ses performances au fil des époques, avec un score F1 de validation qui atteint environ **0.69** dès la 6^e époque, avant de se stabiliser. La différence de performance entre entraînement et validation reste modérée, témoignant d'une meilleure capacité de généralisation par rapport à CamemBERT sur ce corpus enrichi.

8.3.3 Analyse

Ces résultats ont mis en évidence l'intérêt d'utiliser des modèles multilingues plus robustes comme XLM-RoBERTa. Toutefois, les deux modèles montrent des signes de plafonnement en validation, ce qui a motivé l'ajout de données synthétiques générées par des LLM, dans l'objectif d'accroître la diversité des exemples et d'améliorer la robustesse globale du système.

8.4 Sur un jeu de données fusionné (français + traduit + généré)

8.4.1 Optimisation des hyperparamètres et choix du taux d'apprentissage

L'ensemble des expériences repose sur un fine-tuning du modèle **CamemBERT** sur un corpus multilingue d'environ **55 950 phrases** (dont ~5000 originales en français et ~50000 issues de traductions automatiques via **deep-translator** et 950 **générées par LLM**). Le jeu couvre **11 classes** de discours haineux, avec une forte hétérogénéité dans leur répartition. Tous les scripts expérimentaux ont utilisé une configuration commune :

- Tokenisation avec une longueur maximale de 256 tokens
- Batch size compris entre 8 et 16
- Early stopping avec patience de 2 à 5 époques

8.4.2 Exploration initiale des hyperparamètres

Une première série d'expériences (scripts 1 à 5) a été dédiée à l'ajustement des paramètres classiques tels que le taux de dropout et le learning rate. Le tableau 1 résume les performances obtenues :

Script	LR	Dropout	F1 Train	F1 Val	Écart
Initial (1)	2e-5	_	0.80	0.69	0.11
2_hyp (2)	2e-5	0.3	0.78	0.65	0.13
3_hyp (3)	1e-5	0.5	0.75	0.70	0.05
4_hyp (4)	1e-5	0.5	0.73	0.71	0.02
5_hyp (5)	1e-5	0.5	0.72	0.71	0.01

Table 1: Historique des configurations testées avant l'optimisation ciblée du learning rate

Malgré l'amélioration progressive de l'écart train/validation, des limites subsistaient : stagnation des courbes dès la 5^e époque, performances inégales selon les classes, et présence d'un écart train/val récurrent.

8.4.3 Protocole d'optimisation du taux d'apprentissage

Sur recommandation de *Bouzid Hamza* (post-doctorant au GREYC), une nouvelle stratégie expérimentale a été adoptée : **fixer tous les hyperparamètres précédemment optimisés** (dropout, batch size, patience, etc.) pour se concentrer exclusivement sur l'impact du **learning rate** sur la généralisation.

Trois valeurs ont été comparées : **2e-6**, **5e-6** et **1e-5**. L'évolution du *macro-F1 score* sur les ensembles d'entraînement et de validation est illustrée dans la figure 8 :

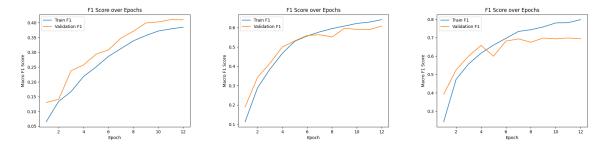


Figure 8: Évolution du macro-F1 score en fonction des époques pour trois valeurs de taux d'apprentissage

8.4.4 Analyse comparative des taux testés

- 2e-6 : taux trop faible, entraînant un apprentissage extrêmement lent. Le F1 d'entraînement reste < 0,42 même après 12 époques, avec une validation qui plafonne à 0,41. Le modèle montre un sous-apprentissage marqué.
- 5e-6 : progression plus rapide, avec un F1 d'entraînement à 0,66, mais le score de validation reste bloqué sous 0,61 dès la 7^e époque. L'écart train/val s'accroît, indiquant un retour du sur-apprentissage.
- 1e-5 : le meilleur compromis observé. F1 d'entraînement atteignant 0,80, avec un score de validation stable autour de 0,69. L'écart train/val est minimal, les courbes sont parallèles, signe d'une bonne convergence et d'une capacité de généralisation équilibrée.

8.4.5 Conclusion

Ces résultats démontrent que le taux **1e-5** est le plus adapté au fine-tuning du modèle sur ce corpus. Il permet une montée en performance efficace tout en évitant le sur-apprentissage, assurant une stabilité des scores de validation et une généralisation satisfaisante.

9 Analyse qualitative

Afin d'évaluer la robustesse sémantique des modèles fine-tunés, une analyse qualitative a été réalisée sur un échantillon de **220 phrases générées manuellement**, réparties

équitablement sur les **11 catégories de discours haineux**. Chaque phrase a été annotée à la main selon le format multi-label du jeu de données.

Deux variantes de corpus ont été utilisées pour entraîner les modèles :

- les données originales en français, fusionnées avec des données **traduites** (EN \rightarrow FR),
- les données précédentes enrichies avec des phrases **générées** artificiellement.

9.1 Résultats avec CamemBERT

9.1.1 Modèle entrainé sur les données FR + Données Traduites

Ce modèle montre une bonne performance générale, notamment sur les catégories fréquentes comme *genre* et *religieux*. La matrice de confusion (Figure 9) illustre une bonne précision globale mais une sensibilité réduite (rappel) sur les classes minoritaires.

- Corrélations: La carte de corrélation (Figure 10) montre des associations modérées entre hateSpeech et genre, religieux ou LGBTQ+, suggérant que le modèle est sensible aux discours haineux à plusieurs facettes.
- $-\mathbf{P}/\mathbf{R}$: Le graphe de précision-rappel (Figure 15) indique une bonne précision sur certaines classes (handicap, genre), mais une faible capacité de rappel sur les classes sensibles comme menace ou violence.

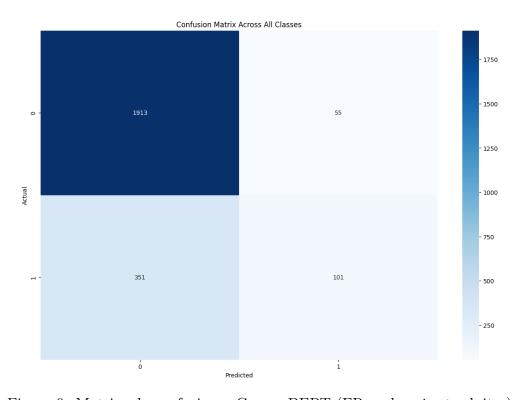


Figure 9: Matrice de confusion – CamemBERT (FR + données traduites)

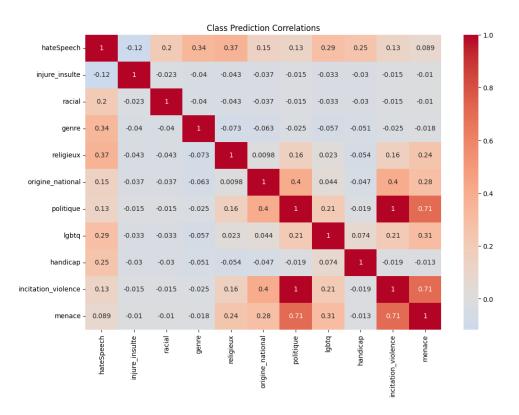


Figure 10: Corrélations de classes – CamemBERT (FR + données traduites)

9.1.2 Modèle entrainé sur les données FR + Données Traduites + Générées

Avec l'ajout de données synthétiques, la performance s'améliore sensiblement, en particulier pour des classes comme menace, LGBTQ+ ou politique.

- Corrélations : Le graphe (Figure 12) met en évidence des liens renforcés entre certaines catégories, notamment *politique*, *violence* et *menace*, ce qui montre que le modèle apprend mieux les cooccurrences.
- **P**/**R** : L'amélioration est manifeste sur la couverture des classes rares (cf. Figure 16) bien que la précision reste parfois à améliorer.
- Confusion: Le modèle a tendance à moins confondre les phrases neutres et haineuses qu'auparavant (Figure 11).

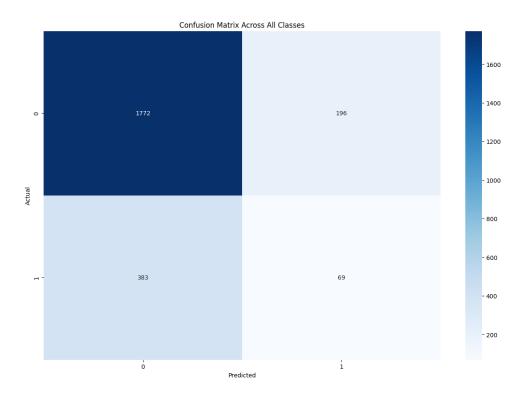


Figure 11: Matrice de confusion – CamemBERT (FR + traduites + générées)

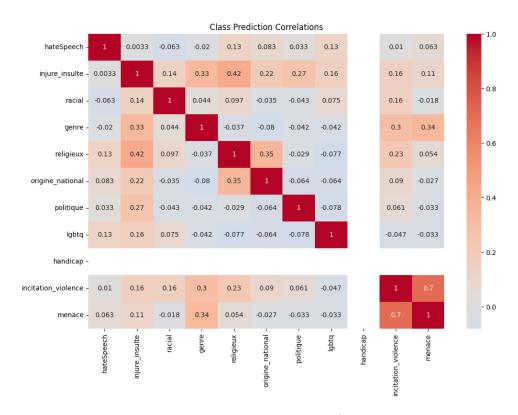


Figure 12: Corrélations de classes – CamemBERT (FR + traduites + générées)

9.2 Résultats avec XLM-Roberta

9.2.1 Modèle entrainé sur les données FR/EN + Données Traduites

Ce modèle multilingue généralise bien, surtout sur les structures syntaxiques complexes. Toutefois, il sous-active certaines classes comme *racial* ou *violence* (Figure 13).

- Corrélations : Légère confusion entre genre et LGBTQ+, comme le montre Figure 14.
- $-\mathbf{P}/\mathbf{R}$: Bonne précision sur *genre* et *religieux*, mais faible rappel pour *menace* et *violence* (Figure 17).

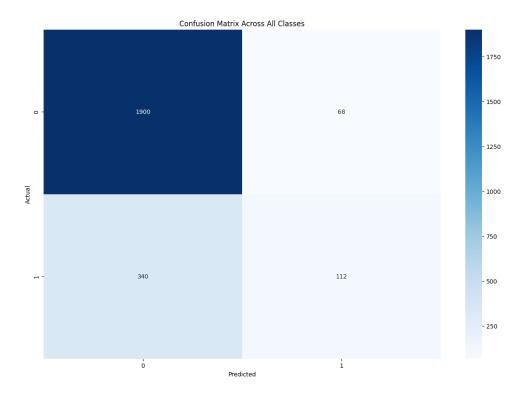


Figure 13: Matrice de confusion – XLM-R (FR/EN + données traduites)



Figure 14: Corrélations de classes – XLM-R (FR/EN + données traduites)

9.2.2 Modèle entrainé sur les données FR/EN + Traduites + Générées

Il s'agit du modèle le plus performant de cette analyse. Il reconnaît mieux les formulations indirectes et les attaques implicites. Les classes auparavant négligées comme handicap et menace sont mieux couvertes, bien que des confusions subsistent entre politique et origine_national.

- Corrélations : Des corrélations fortes apparaissent entre violence, menace, politique et religieux, ce qui traduit une meilleure compréhension contextuelle.
- $-\mathbf{P}/\mathbf{R}$: Très bon équilibre entre précision et rappel, en particulier pour *racial*, *insulte*, *genre*.

9.3 Résumé global

L'ajout de données générées améliore significativement la performance des modèles, en particulier pour les classes rares ou ambiguës. XLM-RoBERTa, enrichi par les données synthétiques, offre la meilleure robustesse globale. CamemBERT montre également des gains, mais reste légèrement moins performant dans les cas complexes.

10 Conclusion

Ce stage au GREYC a été une expérience extrêmement formatrice à l'intersection du TAL et de l'archivistique. J'ai développé une expertise pratique en construisant une pipeline complète de détection de discours haineux - de la collecte et l'augmentation de données multilingues au fine-tuning de modèles transformer (CamemBERT, XLM-RoBERTa).

J'ai confronté les défis réels du ML : déséquilibre des classes, surapprentissage, et limitations éthiques de la génération de données sensibles. L'analyse qualitative a confirmé l'apport des données synthétiques pour les catégories rares.

Techniquement, j'ai acquis une maîtrise approfondie des transformers et de l'optimisation hyperparamétrique. Humainement, j'ai appris à allier performance algorithmique et responsabilité éthique, notamment grâce aux échanges avec l'IMEC. Cette immersion en recherche appliquée a renforcé mon aptitude à mener un projet complexe de A à Z et affiné ma vision des applications sociétales de l'IA.

11 Annexe

Résultats des modèles

Classe	Précision	Rappel	F1-score	Support
hateSpeech	0.78	0.71	0.74	2255
injure_insulte	0.88	0.94	0.91	5200
racial	0.64	0.57	0.60	501
genre	0.83	0.71	0.76	383
religieux	0.92	0.79	0.85	677
origine_national	0.77	0.81	0.79	1205
politique	0.65	0.65	0.65	112
lgbtq	0.80	0.79	0.79	259
handicap	0.81	0.64	0.72	457
incitation_violence	0.50	0.35	0.41	37
menace	0.38	0.38	0.38	8
Micro moy.	0.83	0.82	0.83	11 094
Macro moy.	0.72	0.67	0.69	11094
Pondérée moy.	0.83	0.82	0.83	11094
Échantillons mov.	0.59	0.59	0.59	11094

Table 2: Résultats obtenus par CamemBERT sur l'ensemble de test (1^{er} exécution)

Classe	Précision	Rappel	F1-score	Support
hateSpeech	0.78	0.71	0.74	2255
$injure_insulte$	0.88	0.94	0.91	5200
racial	0.64	0.57	0.60	501
genre	0.83	0.71	0.76	383
religieux	0.92	0.79	0.85	677
$origine_national$	0.77	0.81	0.79	1205
politique	0.65	0.65	0.65	112
lgbtq	0.80	0.79	0.79	259
handicap	0.81	0.64	0.72	457
incitation_violence	0.50	0.35	0.41	37
menace	0.38	0.38	0.38	8
Micro moy.	0.83	0.82	0.83	11 094
Macro moy.	0.72	0.67	0.69	11094
Pondérée moy.	0.83	0.82	0.83	11094
Échantillons moy.	0.59	0.59	0.59	11 094

Table 3: Résultats obtenus par XLM-RoBERTa sur l'ensemble de test $(2^{\text{ème}}$ exécution)

Précision vs rappel

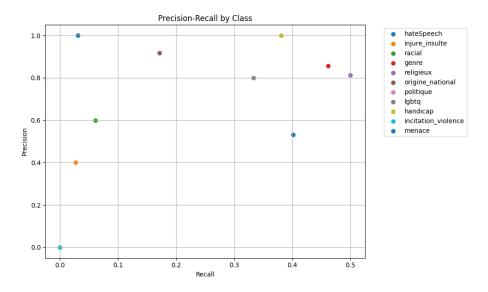


Figure 15: Précision vs rappel – CamemBERT (FR + données traduites)

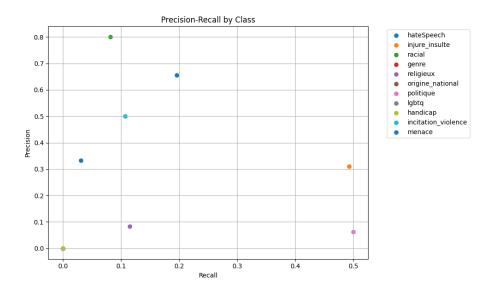


Figure 16: Précision vs rappel – CamemBERT (FR + traduites + générées)

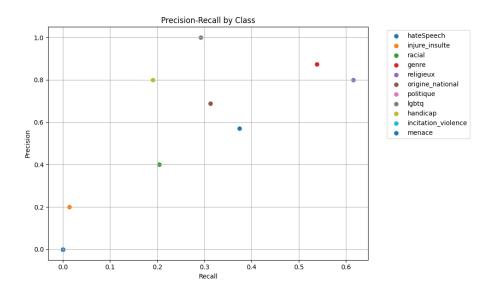


Figure 17: Précision vs rappel – XLM-R (FR/EN + données traduites)

Bibliography

- [1] Guerini, M., & Staiano, J. (2020). CONAN COunter NArratives through Nichesourcing: a Multilingual Dataset of Responses to Fight Online Hate Speech. GitHub repository. https://github.com/marcoguerini/CONAN
- [2] Ouazzani, N. (2022). MLMA_hate_speech: Multilingual and Multi-Aspect Hate Speech Analysis. Hugging Face Datasets. https://huggingface.co/datasets/nedjmaou/MLMA_hate_speech
- [3] Davidson, T., Warmsley, D., Macy, M., & Weber, I. (2017). Automated Hate Speech Detection and the Problem of Offensive Language. Proceedings of the 11th International AAAI Conference on Web and Social Media (ICWSM). https://github.com/t-davidson/hate-speech-and-offensive-language
- [4] Vanetik, N. (2021). FTR-dataset: A French Twitter Racism Dataset. GitHub repository. https://github.com/NataliaVanetik/FTR-dataset
- [5] Vicomtech. (2019). Hate speech dataset from a white supremacist forum. GitHub repository. https://github.com/Vicomtech/hate-speech-dataset
- [6] Mollas, I., Chrysopoulou, Z., Karlos, S., & Tsoumakas, G. (2022). ETHOS: a Multi-label Hate Speech Detection Dataset. arXiv preprint arXiv:2203.03954. https://github.com/intelligence-csd-auth-gr/Ethos-Hate-Speech-Dataset
- [7] Rault, T., & Sochon, P. (2024). Classification des propos injurieux, racistes et homophobes. Travail étudiant, Université de Caen Normandie.
- [8] Martin, L., et al. (2020). CamemBERT: a Tasty French Language Model. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL).
- [9] Conneau, A., et al. (2020). Unsupervised Cross-lingual Representation Learning at Scale. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL).
- [10] Devlin, J., et al. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Proceedings of NAACL-HLT.