

---

# Extraction d'Entités Nommées au format BIO avec spaCy dans différentes versions du texte

Weronika BIEGOWSKA 21211975

[weronika.biegowska@etu.sorbonne-universite.fr](mailto:weronika.biegowska@etu.sorbonne-universite.fr)

Date: 04/04/2025

# Problématique

- La **reconnaissance des entités nommées** (REN) et son importance dans le traitement automatique du langage (TAL).
- La **reconnaissance optique de caractères** (OCR):
  - Karken
  - Tesseract

# Constitution des Corpus

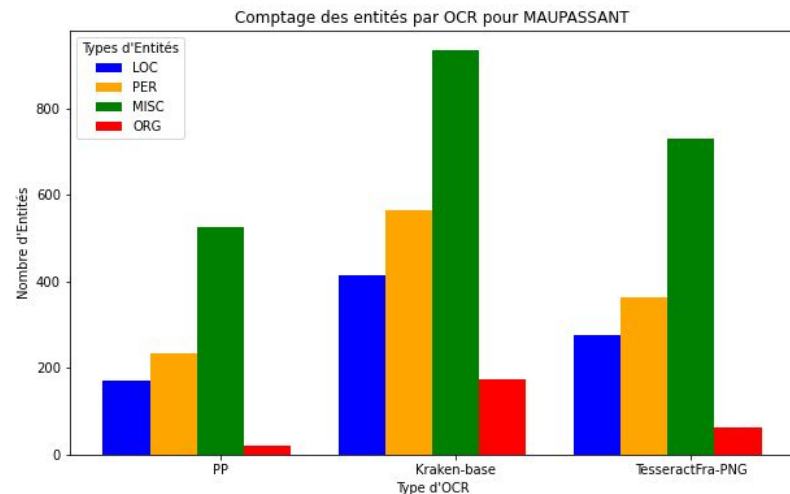
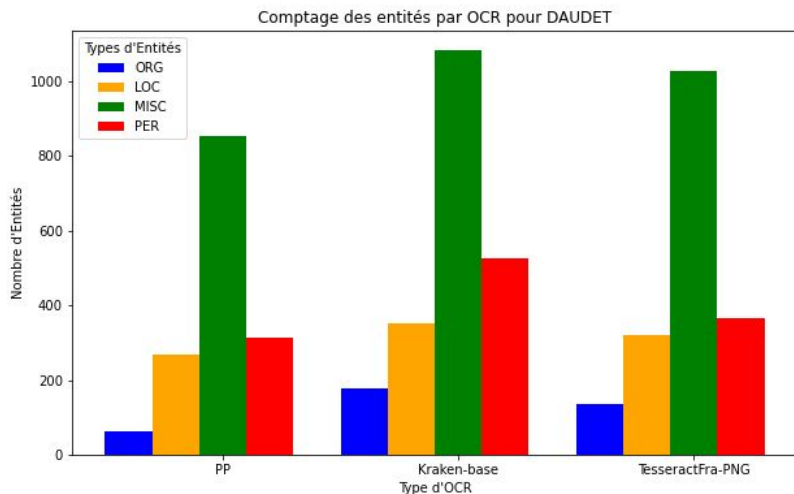
Tab.1: Corpus			
Auteur	Nom du fichier de référence	Nom du fichier de OCR: Kraken	Nom du fichier de OCR: Tesseract
DASH	DASH_chateau-de-Pinon-V1_PP	DASH_chateau-de-Pinon-V1_Kraken-base	DASH_chateau-de-Pinon-V1_TesseractFra-PNG
DAUDET	DAUDET_petit-chose_PP	DAUDET_petit-chose_Kraken-base	DAUDET_petit-chose_TesseractFra-PNG
FLAUBERT	FLAUBERT_education-sentimentale_PP	FLAUBERT_education-sentimentale_Kraken-base	FLAUBERT_education-sentimentale_TesseractFra-PNG
MAUPASSANT	MAUPASSANT_une-vie_PP	MAUPASSANT_une-vie_Kraken-base	MAUPASSANT_une-vie_TesseractFra-PNG
NOAILLES	NOAILLES_la-nouvelle-esperance_PP	NOAILLES_la-nouvelle-esperance_Kraken-base	NOAILLES_la-nouvelle-esperance_TesseractFra-PNG

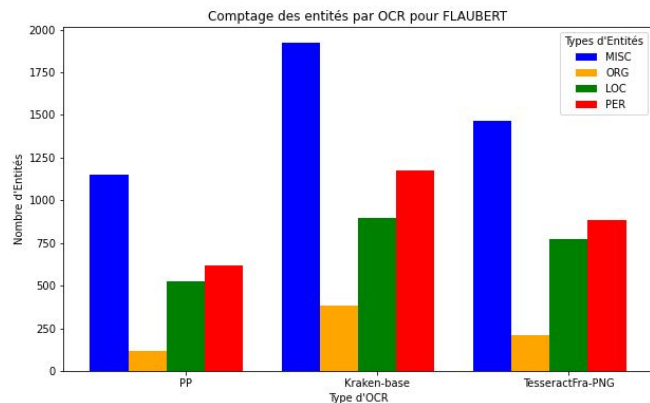
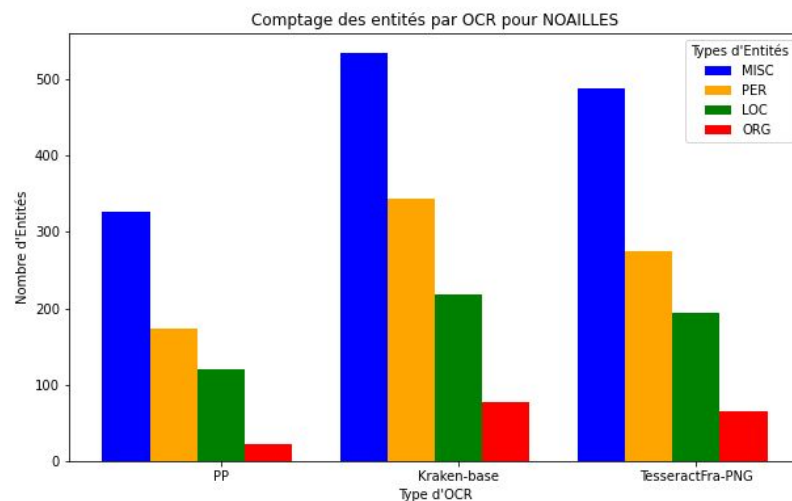
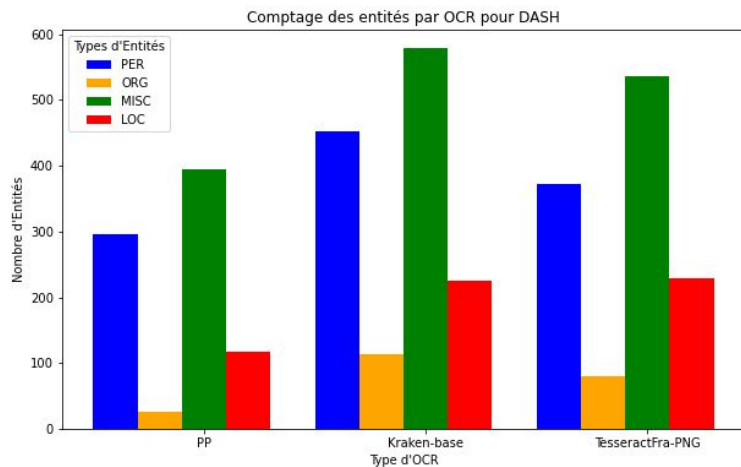
OCR utilisés : **Tesseract** (open-source, généraliste) et **Kraken** (spécialisé pour les documents historiques)

# Méthodes d'analyse

- Utilisation de **spaCy** pour la reconnaissance d'entités.
- Format **BIO** pour l'annotation des entités.
- Reconnaissance d'entités et **POS tagging** (étiquetage de la partie du discours).
- Récupération des tokens avec le label "**PROPN**" (Proper Noun).
- Analyse des erreurs et évaluation de la performance à l'aide de mesures comme la précision, le rappel et le score F1.
- Visualisation

# La proportion d'entités pour chaque label sémantique selon les différentes versions des textes.

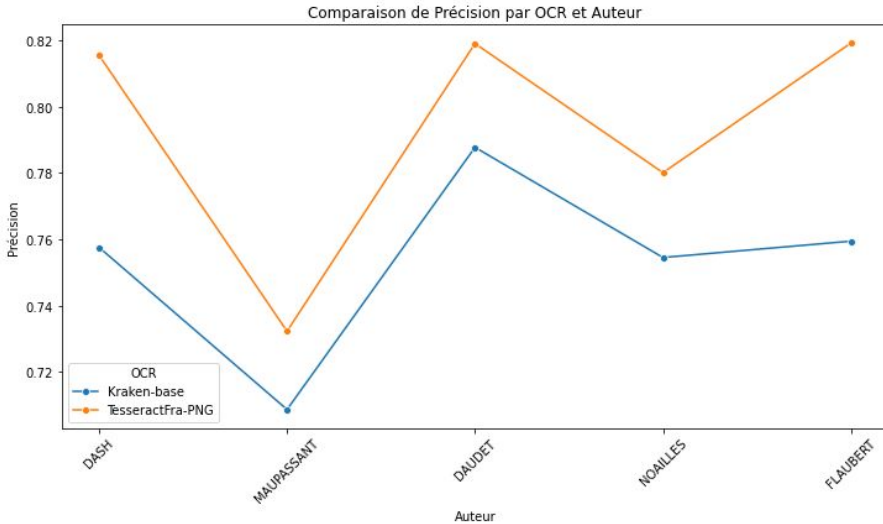




## Résultats clés :

- Plus d'entités nommées ont été retrouvées dans les textes OCR que dans la version de référence.
- Le label sémantique "MISC" est le plus fréquent dans tous les textes et toutes les versions.

# Précision



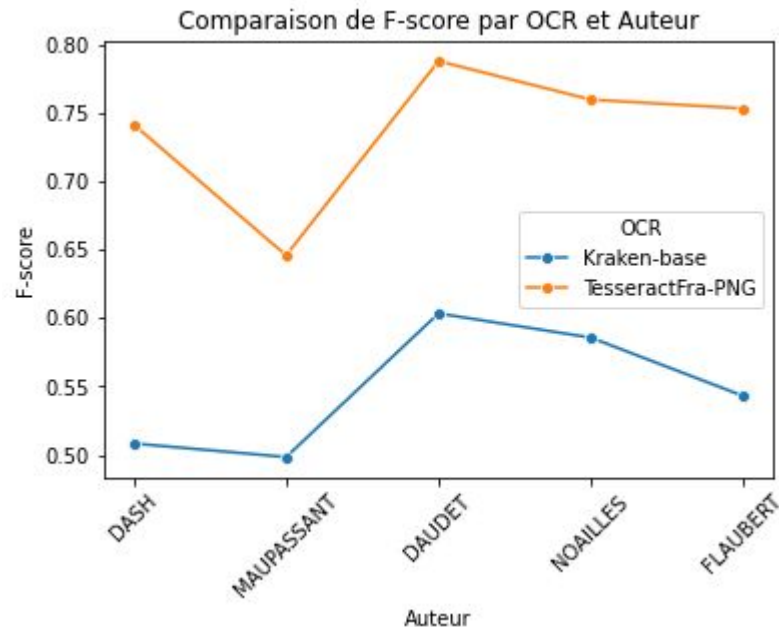
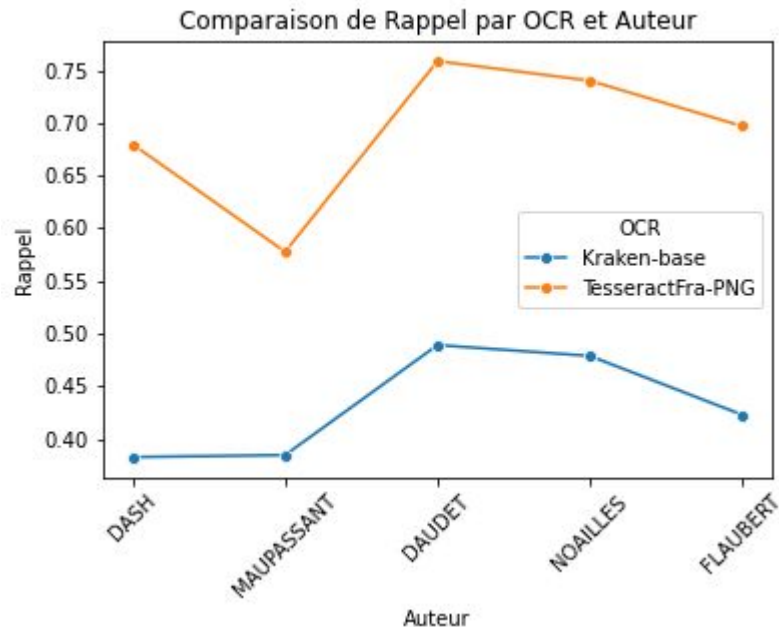
Tab 2: Les résultats des comparaison

Auteur	Référence vs Kraken			Référence vs Tesseract		
	Précision	Rappel	F-score	Précision	Rappel	F-score
DASH	0.75	0.37	0.5	0.81	0.67	0.74
DAUDET	0.78	0.48	0.59	0.82	0.75	0.78
FLAUBERT	0.75	0.41	0.53	0.81	0.69	0.74
MAUPASSANT	0.7	0.37	0.49	0.73	0.57	0.64
NOAILLES	0.75	0.47	0.57	0.77	0.73	0.75

La précision mesure la proportion d'entités identifiées par le modèle qui sont correctes.

**Formule** :  $\text{Précision} = \frac{\text{Vrai Positifs}}{\text{Vrai Positifs} + \text{Faux Positifs}}$

# Rappel et F-score

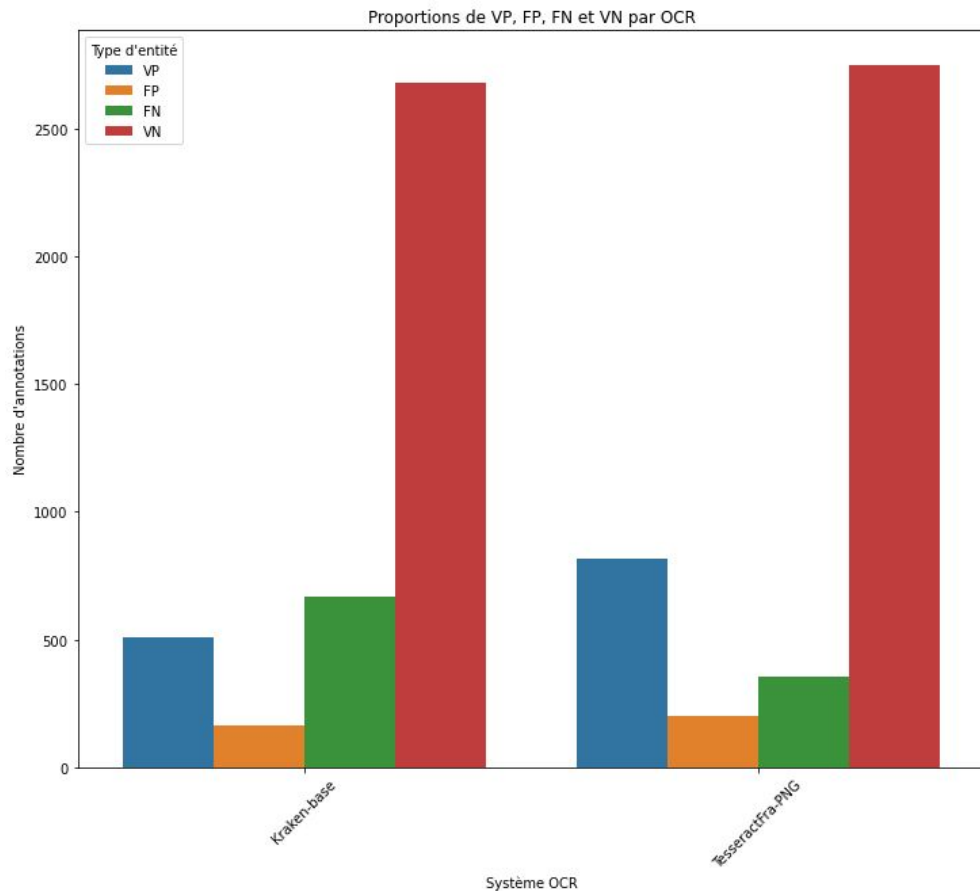


Le **rappel** mesure la proportion des entités correctes qui ont été effectivement trouvées par le modèle.  $\text{Rappel} = \text{Vrai Positifs} / (\text{Vrai Positifs} + \text{Faux Négatifs})$

Le **score F1** est la moyenne harmonique entre la précision et le rappel. Il donne une mesure équilibrée des deux.



# Proportions de VP, FP, VN, FN



# Les intersections entre les sorties de REN et les sorties de POS tagging.

Diagramme de Venn - Kraken-base - DAUDET



Diagramme de Venn - TesseractFra-PNG - DAUDET

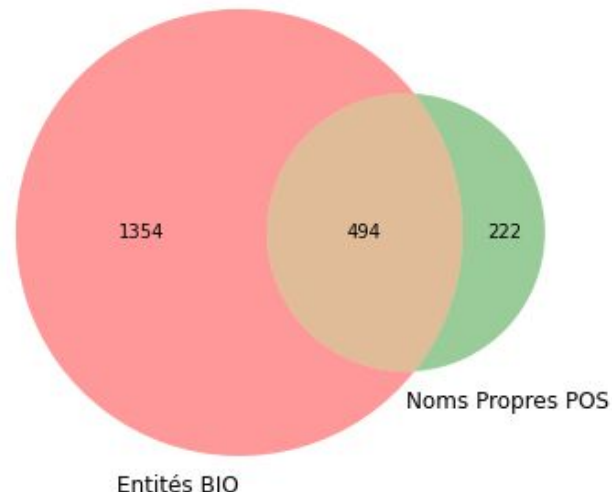


Diagramme de Venn - Kraken-base - MAUPASSANT

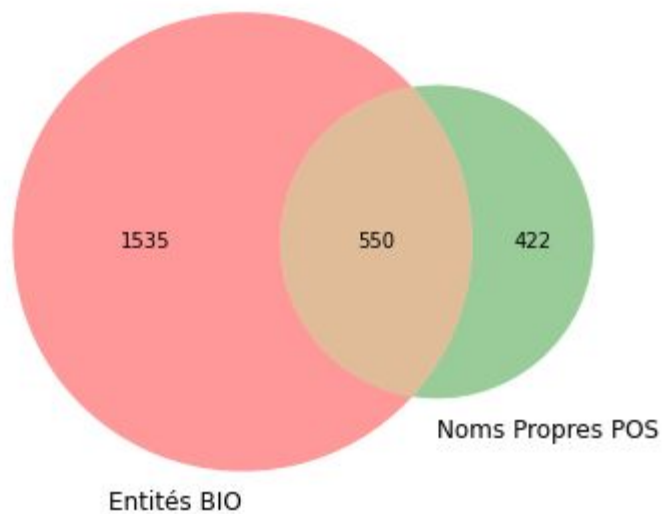


Diagramme de Venn - TesseractFra-PNG - MAUPASSANT

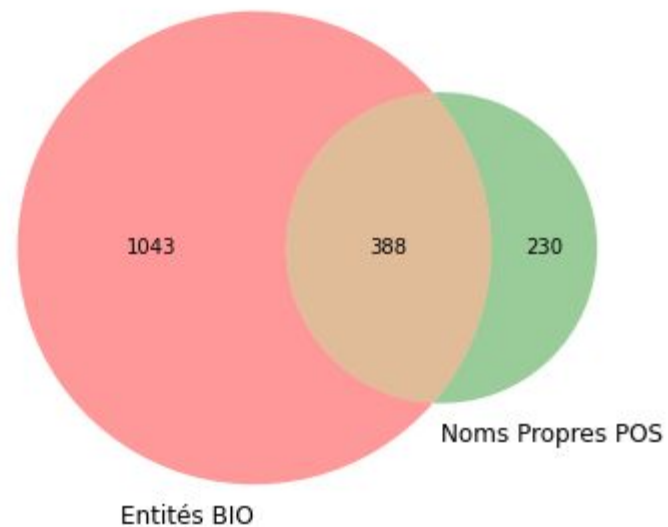


Diagramme de Venn - Kraken-base - DASH

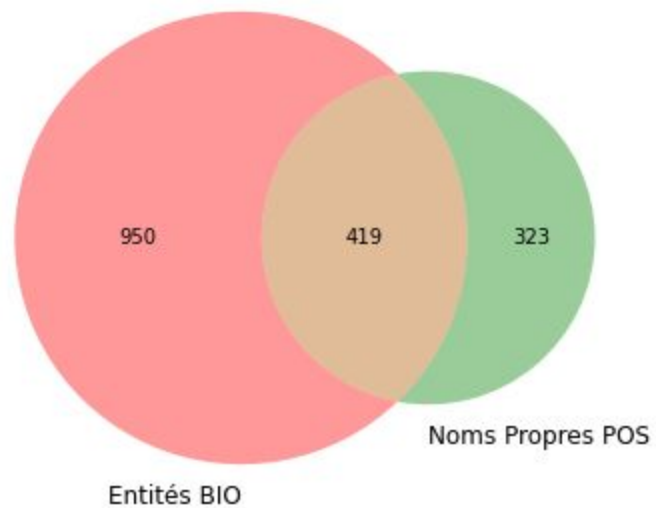


Diagramme de Venn - TesseractFra-PNG - DASH

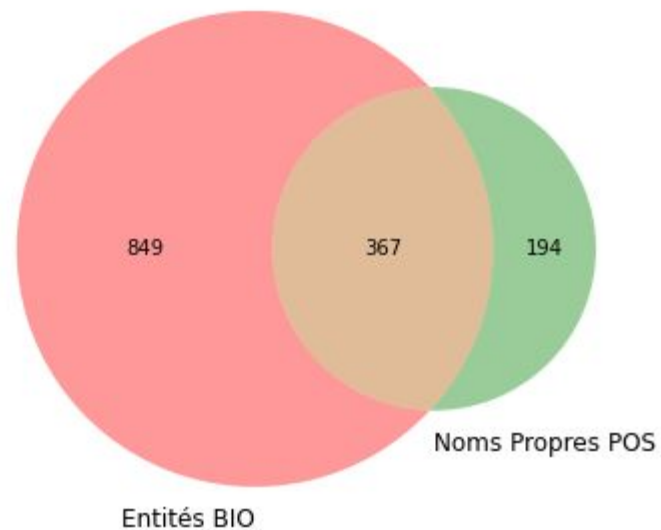


Diagramme de Venn - Kraken-base - NOAILLES

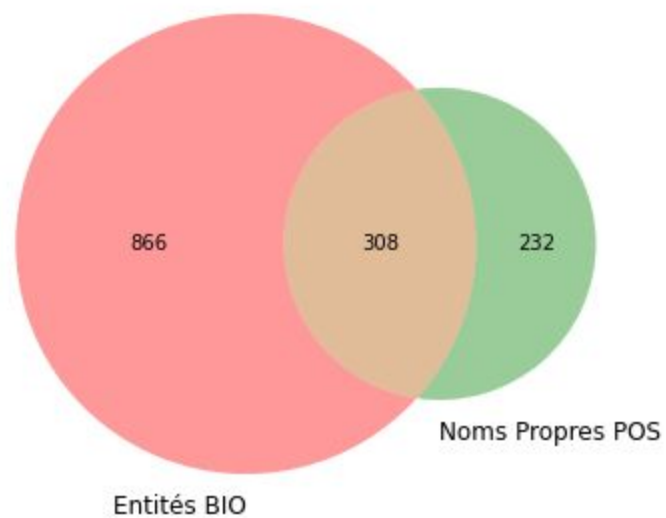


Diagramme de Venn - TesseractFra-PNG - NOAILLES

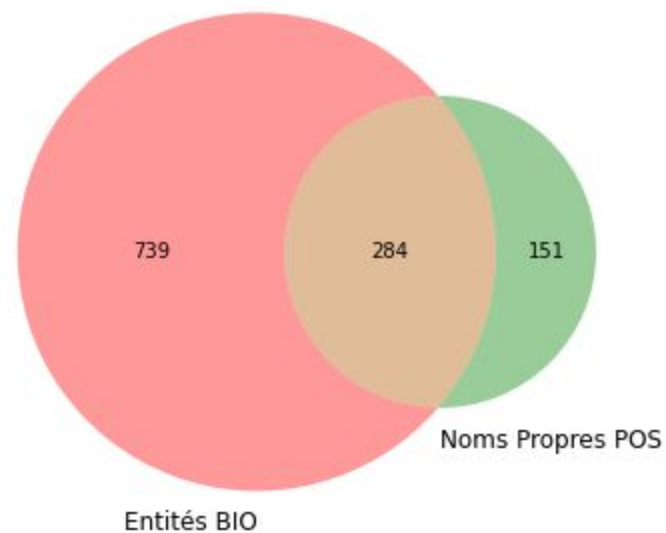


Diagramme de Venn - Kraken-base - FLAUBERT

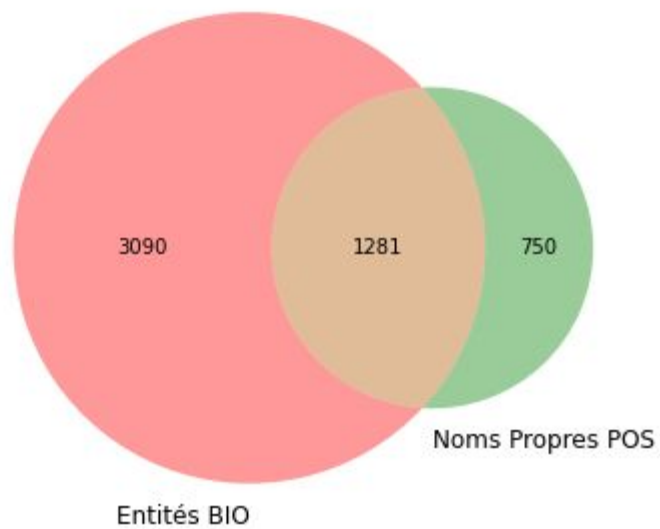
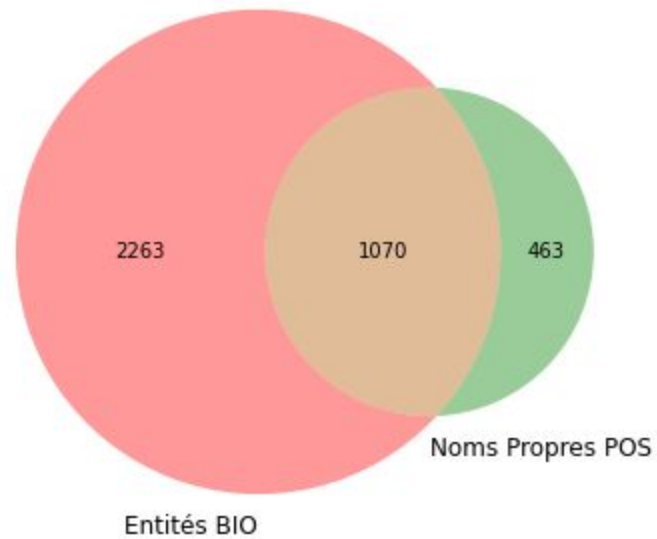


Diagramme de Venn - TesseractFra-PNG - FLAUBERT



# Conclusions

- Cette étude a exploré l'extraction d'entités nommées au format BIO avec spaCy sur différentes versions d'un texte, y compris des versions OCR.
- Les résultats ont montré que la qualité du texte influence fortement la précision du REN, les erreurs OCR perturbant la détection des entités.
- L'analyse via POS tagging a mis en évidence certaines confusions entre entités nommées et noms propres.

# Perspectives

- Entraîner un **modèle spécifique** sur des données bruitées:
- Ajouter une étape de **prétraitement** (correction OCR): une sorte de "nettoyage" du texte qui réduit les fautes typographiques.
- Utiliser des **modèles plus puissants** (type BERT) : peuvent être plus robustes et plus précis que les modèles classiques de spaCy.
- Ce travail souligne les défis du REN sur des textes de qualité variable et ouvre des perspectives pour des approches plus robustes.



# Récapitulation

L'objectif: Cette étude a examiné l'extraction d'entités nommées (REN) avec spaCy sur des textes originaux et des versions issues de la reconnaissance optique de caractères (OCR).

## Principaux résultats:

- La qualité du texte est un facteur déterminant de la précision de la REN.
- Les erreurs introduites par l'OCR ont un impact négatif sur l'identification des entités nommées.
- Tesseract a démontré une meilleure performance globale comparé à Kraken.
- L'analyse morphosyntaxique (POS tagging) a révélé des chevauchements et des distinctions entre les entités nommées et les noms propres.

Des stratégies d'amélioration, telles que l'entraînement sur des données OCR bruitées, peuvent améliorer la robustesse des modèles.