



# La Détection de Personnages dans la Littérature

Quelle Place pour les Modèles de Langue ?

---

Jean Barré, Mathieu Dehouck

Journées MATE-SHS, 23-24 Mai 2024

**Lattice** - UMR 8094 (CNRS, ENS-PSL, Sorbonne Nouvelle)

## Deux approches du personnage dans le contexte littéraire

- Référentielle - Le personnage comme représentation - projection psychologique
- Formelle - La fonction personnage dans le récit, cf l' "être de papier" (Barthes, 1970)

## L'espace personnage dans le roman

- L'espace personnage : « la rencontre entre une personnalité humaine individuelle et une position déterminée dans l'ensemble du récit » (Woloch, 2003)
- Distribution de l'attention accordée à un personnage au fil du récit

## Opérationnalisation de l'espace personnage

- Personnage : au delà de l'entité nommée
- Prise en compte de l'agentivité
- Détection de la coréférence - ensemble des mentions du personnage

Quel importance des personnages secondaires sur le récit ?

## Un algorithme de TAL pour la littérature ?

- Partie Française d'un projet multilingue
- Annotation sur un corpus de romans du XIX<sup>e</sup> et XX<sup>e</sup> siècle
- Détection d'entités (PER, FAC, TIME, ORG, LOC)
- Clustering d'entités
- Résolution de la coréférence pour chaque cluster avec l'aide de CamemBERT



[#37] T131 J 'avais amené T1482 cette jeune femme au bal de T615 madame de Lanty . Comme T1482 elle venait pour la première fois dans T356 cette maison , T131 je T1482 lui pardonnai T1482 son rire étouffé ; mais T131 je T1482 lui fis vivement T131 je ne sais quel signe impérieux qui T1482 la rendit tout interdite et T1482 lui donna du respect pour T1482 son voisin. T1482 Elle s'assit près de T131 moi . T992 Le vieillard ne

**Figure 1** : Chaînes de coréférence dans Sarrasine

## Détection des personnages :

	précision	rappel	$F_1$
PER	85.0	92.1	88,4

## Résolution de la coréférence :

Métriques	$F_1$	
$MUC$	88,0	Average 76.4
$B^3$	69,2	
$CEAF_e$	71.8	

## Manon Lescaut : Un récit test pour Booknlp

- Diversité des personnages (nombreux personnages secondaires)
- Discontinuité des personnages
- Pluralité des dénominations des personnages

### Named characters

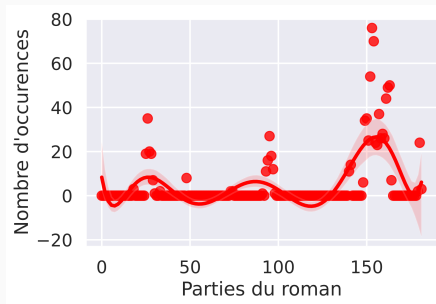
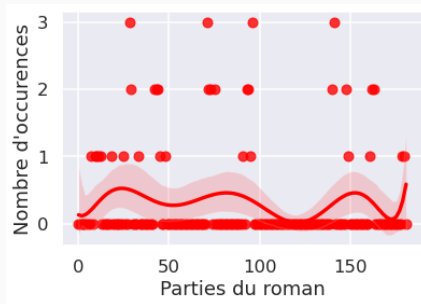
1436 Manon (256)/Lescaut (47)/ma chère Manon (5)/mademoiselle Manon (3)/Manon Lescaut (3)/Lescaut et moi (2)/Lescaut et Manon (2)/l' infidèle Manon (1)/belle Manon (1)/Chère Manon (1)/ta Manon (1)/la pauvre Manon (1)/La pauvre Manon (1)/la perfide Manon (1)/Inconstante Manon (1)/la malheureuse Manon (1)/chère Manon (1)

**Figure 2 :** Entités associées au nom propre « Manon »

## L'espace personnage des seconds rôles

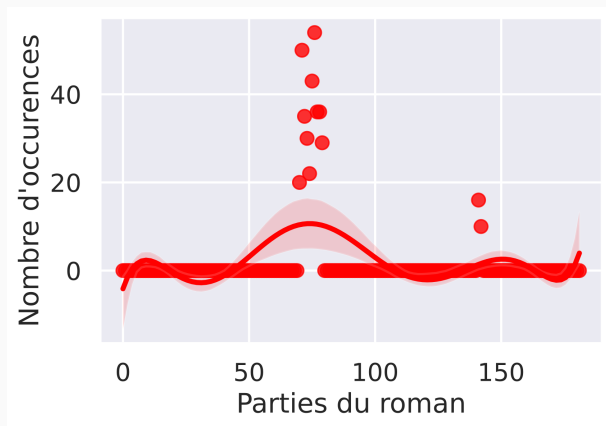
- Inférences de fr-BookNLP sur tout le récit
- Fenêtre roulante de 1 000 tokens - Signal de l'espace personnage
- Intermittence des personnages secondaires ?
- Améliorer notre compréhension de la structure narrative

# Tiberge dans Manon Lescaut



**Figure 3 :** Occurrences vs Coréférence de Tiberge au fil du roman

## Des Personnages sans nom : le lieutenant de police



**Figure 4 :** Occurrences de la coréférence du lieutenant de police au fil du roman

- Apporter un nouveau point de vue sur les études littéraires
- Analyse à grande échelle des personnages
- Étude de leur caractérisation -> les mots qui les décrivent, les actions qu'ils entreprennent
- Construire à la Propp ou à la Greimas une typologie des personnages
- Différenciation en fonction de leur genre
- Promesses et limites des LLMs : Quid des modèles génératifs ?

## Quelques considérations

- Accès local ou par une API ?
- Open source ? Besoins en mémoire et en calcul ?
- Utilisation payante ou gratuite ?
- Adaptable à nos données (fine-tunable) ou pas ?
- Résultats reproductibles ? Déterministe, stochastique, mis à jour ?
- Accès : à une partie de la sortie ? à toute la sortie ? aux couches intermédiaires ?
- Que sait-on des données d'apprentissage ?
- Résultats : état de l'art ? utilisables ? passables ?



# Comment tester son LLM ?

*Speak, Memory : An Archaeology of Books Known to ChatGPT/GPT-4.*  
Chang, Cramer, Soni & Bamman.

Impact des données d'apprentissage sur l'évaluation des LLMs.

## Exemples d'expériences

- Prédire un nom propre dans un court passage ;
- Prédire la date de première parution.

## Résultats : ChatGPT connaît très bien

- Les œuvres libres de droits : Alice au pays des merveilles (L. Carroll), Les Aventures de Sherlock Holmes (A. Conan Doyle), Frankenstein (M. Shelley), Orgueil et Préjugés (J. Austen)...
- Les œuvres bien identifiées en ligne : Harry Potter à l'école des sorciers (J. K. Rowling), La Communauté de l'Anneau (J. R. R. Tolkien)...

## Expérience de coin de table

1. Choisir deux objets/entités de la même catégorie, l'un très connu, l'autre beaucoup moins;
2. Demander à un modèle de diffusion d'images de les dessiner.

## Exemple

- Une carte des États-Unis d'Amérique;
- Une carte du Bhoutan.
- <https://huggingface.co/spaces/ByteDance/SDXL-Lightning>

# Les États-Unis d'Amérique (1)



# Les États-Unis d'Amérique (2)



Figure 6 : Générée par Stable Diffusion

# Le Bhoutan (pour référence)



**Figure 7 :** Prise sur Wikidata

# Le Bhoutan (1)



Figure 8 : Générée par Stable Diffusion

# Le Bhoutan (2)



- Au delà des biais et de la contamination;
- Fort impact de l'environnement culturel d'entraînement;
- Tester ses modèles sur des objets culturellement bien identifiés mais aussi sur des objets mal identifiés.



### **Projet à Berkeley :**

- <https://github.com/dbamman/litbank>
- <https://github.com/booknlp/booknlp>

### **Partie française développée au Lattice :**

- <https://www.lattice.cnrs.fr/projets/booknlp/>
- <https://github.com/lattice-8094/fr-litbank/>