

# Simulation d'erreurs d'OCR dans les systèmes de TAL pour le traitement de données anachroniques

Baptiste Blouin, Benoit Favre, Jeremy Auguste

27-06-2022





# Les erreurs d'OCR

The Commander-in-Chief of the UNITED STATES forces had declived to furnish arms under his charge, at the requisition of the Governor. The VIGILANOE COMMITTEE Were



(a)

The Commander-in-Chief of the UNITED STATES forces had declived to furnish arms under his charge, at the requisition of the Governor. The VIGILANOE COMMITTEE Were

REFERRING to the remarks in our last issue on the state of affairs in CALIFORNIA, we resume the subject. We gather from various sources, that,



(b)

REFERRING to the remarks in our last issue on the state of affairs in CALIFORNIA, we resume the subject. We gather from various sources, that,

At a quarter past one o'clock, Casey and Corn were brought out upon the platforms. The former was attended by the Rev. Father Gallagher.



(c)

At a quarter past one o'clock, Casey and Corn were brought out upon the platforms. The former was attended by the Rev. Father Gallagher.

- À quel point les erreurs d'OCR impact les tâches de TAL ?
  - ▶ Reconnaissance d'entités nommées [Grover et al., 2008, Packer et al., 2010]
  - ▶ Étiquetage en partie de discours [Lin, 2003, Mieskes and Schmunk, 2019]
  - ▶ Résumé automatique [Jing et al., 2003]
  - ▶ [Linhares Pontes et al., 2019, Boros et al., 2022]
- Est-ce que les modèles clef en main sont utilisables sur ce type de données ? [Rodriguez et al., 2012]
- Est-ce un problème de modèle ou un problème des données d'apprentissage ?

## Annotations

- Un faible quantité de données annotés pour le domaine historique.
  - ▶ Quaero Old Press [Rosset et al., 2012]
  - ▶ HIPE [Ehrmann et al., 2020]
  - ▶ NewsEye [Hamdi et al., 2021]
- Comment tirer parti des ressources contemporaines annotées ?

## Proposition

- Évaluation des nouvelles architectures *Transformers* sur ces difficultés.
- Injection d'erreurs d'OCR dans les données annotées contemporaines.

# Configuration

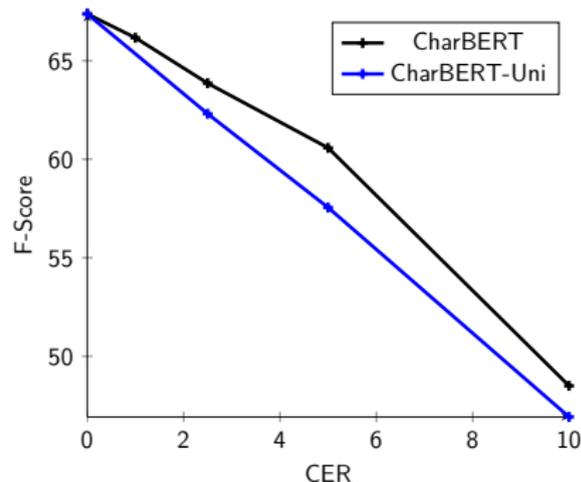
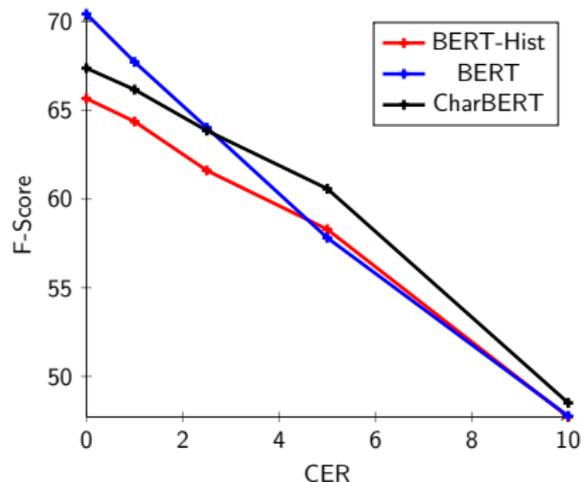
- Dataset
  - ▶ **OntoNote 05** [Weischedel et al., 2013]
  - ▶ **LitBank** [Bamman et al., 2019]
  - ▶ **ACE 05** [Walker et al., 2006]
  - ▶ Le Wall Street Journal du jeu de données Penn Treebank [Marcus et al., 1993]
- Tache
  - ▶ Reconnaissance d'entité nommée
  - ▶ Étiquetage en partie de discours
  - ▶ Extraction d'arguments d'événement
- Modèles
  - ▶ **BERT** [Devlin et al., 2019]
  - ▶ **BERT-Hist** [Hosseini et al., 2021].
    - Microsoft British Library de 1760 à 1900
  - ▶ **CharBERT** [Ma et al., 2020]

# Injection d'erreurs d'OCR

- Modification de la séquence de caractères originale en utilisant la suppression, l'insertion et la substitution [Pruthi et al., 2019].
- Utilisation des probabilités d'erreurs provenant de ICDAR-17 [Nayef et al., 2017].
  - ▶ Presses anglaises, de la British Library, allant de 1744 à 1911.
  - ▶ Permet de contrôler la quantité de CER à injecter.
- Une injection d'erreurs dans les jeux de tests permet d'évaluer la robustesse des modèles dans un scénario d'adaptation de domaine.

# Quantification de l'impact des erreurs d'OCR

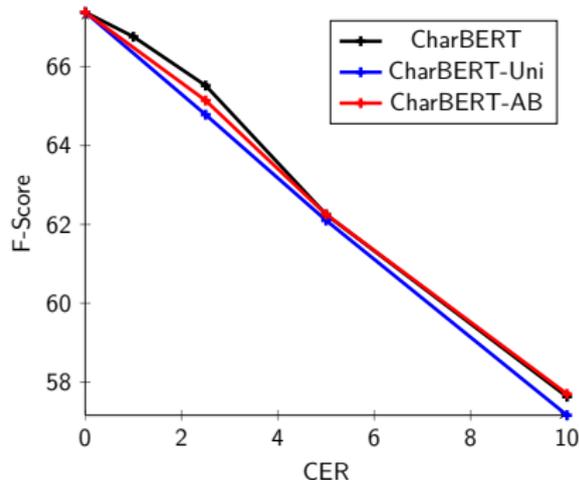
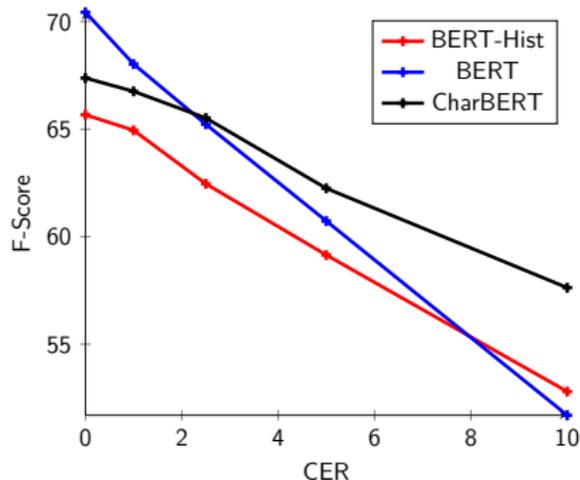
## Reconnaissance d'entités nommées : ACE → LitBank



- Atténuation de l'impact du bruit avec les représentations de mots historiques ou avec l'utilisation de CharBERT.

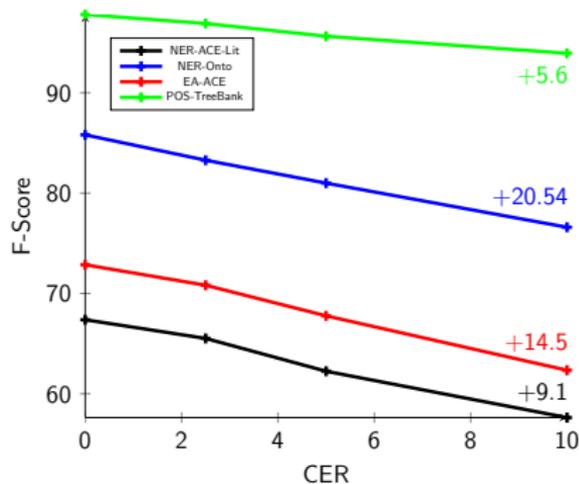
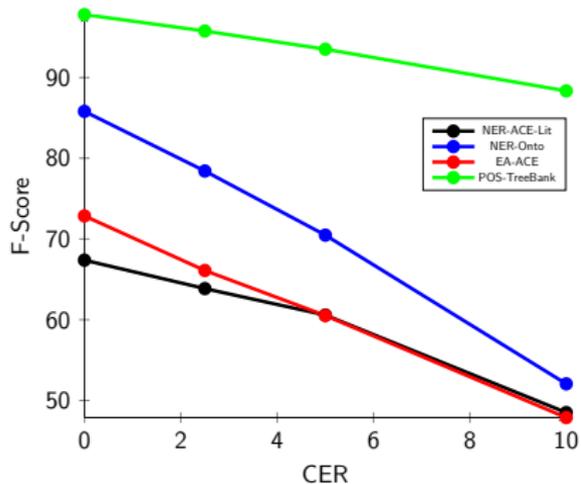
# Réduction de l'impact des erreurs d'OCR

Reconnaissance d'entités nommées : ACE → LitBank



- CharBERT est plus robuste que BERT lors d'un apprentissage sur des données bruitées (+9.1% comparé à +3.94 %).
- La quantité de CER a plus d'impact que la distribution d'erreurs.

## Extraction d'informations



- Scénario d'application dans le même domaine.
- Toutes les tâches sont impactées par les erreurs d'OCR.

- Une quantité de bruit non homogène dans les documents historiques.

### Extraction d'arguments d'événements : ACE → ACE

Train \ Test	0%	2.5%	5%	10%
0%	<b>72.86%</b>	67.12%	63.16%	49.72%
2.5%	68.47%	<b>70.82%</b>	66.65%	56.71%
5%	68.30%	68.79%	<b>67.77%</b>	59.15%
10%	67.36%	68.58%	67.22%	<b>62.38%</b>

- Plus la quantité de bruit est proche de celle des documents d'inférences, plus les résultats sont élevés.
- Un entraînement sur des données bruitées pour le traitement de données non bruitées a moins d'impact que l'inverse.

# Conclusion

- Évaluation de l'impact des erreurs d'OCR sur des tâches d'extraction d'informations.
- Amélioration des systèmes de TAL en injectant du bruit dans les données d'entraînement contemporaine.
- La quantité de bruit à plus d'impact que la provenance des erreurs.
- L'utilisation d'embeddings de caractères permet de réduire cet impact.

## Travaux futurs

- Correction ciblée de ces erreurs d'OCR sur les informations extraites.
- Analyse d'erreurs approfondie afin de mieux identifier les difficultés liées aux erreurs d'OCR.

# Questions ?



# References

-  Bamman, D., Papat, S., and Shen, S. (2019).  
An annotated dataset of literary entities.  
In *Proceedings of the 2019 Conference of the North*, pages 2138–2144,  
Minneapolis, Minnesota. Association for Computational Linguistics.
-  Boros, E., Nguyen, N. K., Lejeune, G., and Doucet, A. (2022).  
Assessing the impact of OCR noise on multilingual event detection over  
digitised documents.  
*International Journal on Digital Libraries*.
-  Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019).  
BERT: Pre-training of deep bidirectional transformers for language  
understanding.  
In *Proceedings of the 2019 Conference of the North American Chapter of the  
Association for Computational Linguistics: Human Language Technologies,  
Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis,  
Minnesota. Association for Computational Linguistics.
-  Ehrmann, M., Romanello, M., Fluckiger, A., and Clemenide, S. (2020).  
Extended Overview of CLEF HIPE 2020: Named Entity Processing on