

# Analyse en dépendances du français avec des plongements contextualisés

Loïc Grobol<sup>1, 2, 3</sup> Benoît Crabbé<sup>1</sup>

(1) LLF, CNRS, Université de Paris, 8, Rue Albert Einstein 75013 Paris, France

(2) Lattice, CNRS, ENS, PSL, Université Sorbonne Nouvelle, 1 Rue Maurice Arnoux, 92120 Montrouge, France

(3) LIFO, ICVL, Université d'Orléans, 45000 Orléans, France

loic.grobol@ens.psl.eu, benoit.crabbe@linguist.univ-paris-diderot.fr

## RÉSUMÉ

---

Cet article présente un analyseur syntaxique en dépendances pour le français qui se compare favorablement à l'état de l'art sur la plupart des corpus de référence. L'analyseur s'appuie sur de riches représentations lexicales issues notamment de BERT et de FASTTEXT. On remarque que les représentations lexicales produites par FLAUBERT ont un caractère auto-suffisant pour réaliser la tâche d'analyse syntaxique de manière optimale.

## ABSTRACT

---

### French dependency parsing with contextualized embeddings

This paper presents a dependency parser for French that compares favorably to the state of the art on several corpora. The parser relies on rich lexical representations from BERT and FASTTEXT. We notice that the lexical representations produced by FLAUBERT are somehow self-sufficient to perform the syntax analysis task in an optimal way.

---

**MOTS-CLÉS** : Analyse syntaxique en dépendances du français, BERT, FastText.

**KEYWORDS** : Dependency Parsing, Parsing French, BERT, FastText.

---

## 1 Introduction

Cet article décrit un modèle d'analyse syntaxique en dépendances couplé à un étiqueteur morphosyntaxique pour le français<sup>1</sup>. Nous étudions plus spécifiquement l'impact de représentations lexicales apprises de manière non supervisée sur de gros volumes de texte telles que FASTTEXT et BERT en complément de plongements lexicaux standards. Nous commençons par expliciter le modèle d'analyse en section 2 avant de présenter les expériences sur les représentations lexicales en section 3.

## 2 Modèle d'analyse

Le modèle d'analyse combine à la fois un étiqueteur morphosyntaxique, un analyseur basé sur l'algorithme de [Dozat & Manning \(2017\)](#) et l'utilisation de riches représentations lexicales ([Bojanowski et al., 2017](#); [Devlin et al., 2019](#)).

---

1. Disponible à <https://github.com/bencrabbe/npdependency>

Dans ce qui suit, on suppose qu'une phrase  $w_1 \dots w_n$  est représentée par une séquence de plongements lexicaux  $\mathbf{X} = \mathbf{x}_1 \dots \mathbf{x}_n$ . En préalable à l'analyse nous traitons les plongements lexicaux à l'aide d'un LSTM pour construire une séquence de représentations contextualisés  $\mathbf{C} = \mathbf{c}_1 \dots \mathbf{c}_n$  :

$$\mathbf{C} = \text{LSTM}(\mathbf{X})$$

L'étiquetage morphosyntaxique est réalisé à l'aide d'un réseau à propagation avant (MLP) et d'une sortie SOFTMAX :

$$\hat{e}_i = \text{SOFTMAX}(\text{MLP}^{\text{POS}}(\mathbf{c}_i))$$

$\hat{e}_i$  est alors la distribution de scores des étiquettes morphosyntaxiques pour le mot  $w_i$ .

Pour l'analyse syntaxique, nous commençons par spécialiser les représentations de  $\mathbf{C}$  de quatre manières différentes :

$$\mathbf{h}_i^{\text{arc-dep}} = \text{MLP}^{\text{arc-dep}}(\mathbf{c}_i)$$

$$\mathbf{h}_i^{\text{arc-gov}} = \text{MLP}^{\text{arc-gov}}(\mathbf{c}_i)$$

$$\mathbf{h}_i^{\text{label-dep}} = \text{MLP}^{\text{label-dep}}(\mathbf{c}_i)$$

$$\mathbf{h}_i^{\text{label-gov}} = \text{MLP}^{\text{label-gov}}(\mathbf{c}_i)$$

où  $\mathbf{h}_i^{\text{arc-gov}}$  (resp.  $\mathbf{h}_i^{\text{arc-dep}}$ ) représente une spécialisation de la représentations de  $w_i$  du mot comme gouverneur (resp. comme dépendant). Ces représentations étant utilisées pour prédire les arcs, l'utilisation de représentations différentes pour les rôles de gouverneur et de dépendant permet d'éviter des effets de symétrie : le score d'un arc entre deux mots ne devrait en effet pas être le même selon que l'on choisisse soit l'un soit l'autre comme gouverneur. On réalise également une spécialisation analogue pour créer des représentations spécifiques pour la tâche d'étiquetage des arcs que l'on note  $\mathbf{h}_i^{\text{label-dep}}$  et  $\mathbf{h}_i^{\text{label-gov}}$ .

La prédiction des arcs et de leurs étiquettes repose sur plusieurs fonctions biaffines de la forme :

$$\text{BIAFF}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{W} \mathbf{y} + \mathbf{U}(\mathbf{x} \oplus \mathbf{y}) + \mathbf{b}$$

On utilise une première fonction biaffine pour attribuer un score  $s_{j \rightarrow i}^{\text{arc}}$  à chaque arc entre un gouverneur  $w_j$  et un dépendant  $w_i$  :

$$s_{j \rightarrow i}^{\text{arc}} = \text{BIAFF}(\mathbf{h}_i^{\text{arc-gov}}, \mathbf{h}_j^{\text{arc-dep}})$$

et autant de fonctions biaffines qu'il y a de types de dépendances  $\ell$  pour attribuer un score à chaque étiquette possible pour ce même arc :

$$s_{j \rightarrow i}^{\ell} = \text{BIAFF}_{\ell}(\mathbf{h}_i^{\text{label-gov}}, \mathbf{h}_j^{\text{label-dep}})$$

Ces fonctions diffèrent les unes des autres, les matrices de paramètres  $\mathbf{W}_{\ell}$  et  $\mathbf{U}_{\ell}$  étant *a priori* distinctes pour chaque  $\ell$ .

**Fonction objectif** La fonction objectif du modèle est composée d'objectifs multiples. Le premier porte sur l'étiquetage morphosyntaxique. On suppose qu'une phrase est un couple  $(W, E)$  composé d'une séquence de mots  $W = w_1 \dots w_n$  et d'une séquence d'étiquettes de référence  $E = e_1 \dots e_n$ . Cette première fonction calcule l'entropie croisée entre la distribution prédite par le modèle  $\hat{e}_i$  pour chaque tag et la distribution ponctuelle  $\mathbf{e}_i = (\delta_{ij})_j$  qui encode  $e_i$  :

$$\mathcal{L}(W, T) = \sum_{i=1}^n H(\mathbf{e}_i, \hat{\mathbf{e}}_i)$$

Le second objectif porte sur la prédiction des arcs, indépendamment de leurs étiquettes. Dans ce cas on suppose que la phrase est un couple  $(W, D)$ , où  $D$  est un ensemble d’arcs de dépendance de référence. Pour chaque arc  $(j \rightarrow i) \in D$  on construit un vecteur  $\mathbf{g}_i = (\delta_{jk})_k$  qui indique la position du gouverneur  $w_j$  du dépendant  $w_i$ . On obtient du modèle une matrice  $\hat{\mathbf{G}} = \hat{\mathbf{g}}_{i,j}$  telle que pour tout  $i$   $(\hat{\mathbf{g}}_{i,j})_j = \text{SOFTMAX}(s_{1 \rightarrow i}^{\text{arc}} \dots s_{n \rightarrow i}^{\text{arc}})$  et la fonction objectif est alors de la forme suivante :

$$\mathcal{L}(W, D) = \sum_{i=1}^n H(\mathbf{g}_i, \hat{\mathbf{g}}_i)$$

Le dernier objectif concerne l’étiquetage des dépendances. Celui-ci suppose un étiquetage de l’ensemble des arcs de référence par une fonction  $L$ . On code l’étiquetage de référence d’un arc  $(j \rightarrow i) \in D$  par une étiquette  $\ell = L(j \rightarrow i)$  par le vecteur  $\mathbf{l}_{j \rightarrow i} = (\delta_{k\ell})_k$ . Pour chaque arc  $(j \rightarrow i) \in D$  on obtient du modèle un vecteur de score  $\hat{\mathbf{l}}_{j \rightarrow i} = \text{SOFTMAX}(s_{j \rightarrow i}^{\ell_1} \dots s_{j \rightarrow i}^{\ell_k})$ . Pour une phrase donnée on calcule l’entropie croisée, c’est-à-dire :

$$\mathcal{L}(W, D, L) = \sum_{(j \rightarrow i) \in D} H(\mathbf{l}_{j \rightarrow i}, \hat{\mathbf{l}}_{j \rightarrow i})$$

Considérons un corpus arboré  $T = ((W_i, E_i, D_i, L_i))_{1 \leq i \leq N}$  annoté en dépendances et étiqueté morphosyntaxiquement. La fonction objectif globale est la somme des trois objectifs définis ci-dessus :

$$\mathcal{L}(T) = \sum_{i=1}^N \mathcal{L}(W_i, E_i) + \mathcal{L}(W_i, D_i) + \mathcal{L}(W_i, D_i, L_i)$$

Le problème d’optimisation consiste alors à minimiser  $\mathcal{L}(T)$ . En pratique, nous utilisons la variante de descente de gradient stochastique appelée Adam (Kingma & Ba, 2014) et retenons le modèle qui minimise la perte sur le jeu de validation au cours d’un nombre  $e$  d’épochs fixé à l’avance.

**Prédiction** Pour prédire un arbre de dépendance à partir d’une phrase  $w_1 \dots w_n$  donnée en entrée, on commence par évaluer  $s_{j \rightarrow i}^{\text{arc}}$  pour  $1 \leq i \leq n$  et  $1 \leq j \leq n$ . La matrice de scores  $\hat{\mathbf{G}}$  est vue comme la matrice de poids d’un graphe pondéré complet pour lequel on calcule une arborescence  $A$  couvrante de poids maximal à l’aide de l’algorithme de Chu-Liu/Edmonds (Chu & Liu, 1965 ; Edmonds, 1967). On assigne alors à chaque arc  $(j \rightarrow i) \in A$  l’étiquette de score maximal :

## 2.1 Représentations lexicales

L’analyseur utilise des représentations lexicales de natures diverses, et la représentation  $\mathbf{x}_i$  d’un mot  $w_i$  dans la phrase est la concaténation de représentations calculés par différents modèles lexicaux :

$$\mathbf{x}_i = \text{FASTTEXT}(w_i) \oplus \text{BERT}(w_i) \oplus \text{CHAR-RNN}(w_i) \oplus \text{LOOKUP}(w_i) \quad (1)$$

La représentation LOOKUP est une représentation vectorielle ne dépendant que de la forme lexicale de  $w_i$  et stockée dans un dictionnaire. La représentation au niveau des caractères (CHAR-RNN) est obtenu par encodage de la séquence des caractères du mot à l’aide d’un bi-LSTM (Hochreiter & Schmidhuber, 1997). Les paramètres de ces deux représentations sont initialisées aléatoirement et appris sur les données du treebank d’entraînement en même temps que les paramètres de l’analyseur.

En revanche, les représentations FASTTEXT (Bojanowski *et al.*, 2017) et BERT (Devlin *et al.*, 2019) sont obtenus à partir de modèles entraînés sur de beaucoup plus gros volumes de données, et

sont seulement ajustés sur le treebank d’entraînement. L’embedding FASTTEXT est la moyenne des embeddings  $\mathbf{x}_s$  de l’ensemble  $S(w_i)$  des sous mots qui composent le mot  $w_i$ . Les embeddings BERT sont calculés à l’aide d’une succession de réseaux Transformer (Vaswani *et al.*, 2017). Précisément, on utilise une séquence de transformations de la forme :

$$\begin{aligned} \mathbf{c}_i^0 &= \text{LOOKUP}(w'_i) \oplus \text{POSITION}_i \\ \mathbf{c}_i^{l+1} &= \text{TRANSFORMER}(\mathbf{c}_i^l; \mathbf{c}^l) \quad (0 \leq l < 12) \\ \text{BERT}(w'_i) &= \frac{1}{12} \sum_l \mathbf{c}_i^l \end{aligned}$$

où POSITION est une famille de vecteurs encodant des positions et où  $w'$  n’est plus une séquence de mots, mais une séquence de *sous-mots* fournie par un segmenteur appris automatiquement. Enfin, comme représentation  $\text{BERT}(w_i)$  du mot  $w_i$ , on choisit la moyenne des embeddings de ses sous-mots.

### 3 Expériences

Nous avons testé l’analyseur sur différents jeux de données, principalement issues du projet Universal Dependencies (Zeman *et al.*, 2020). Les expériences de développement sont réalisées sur le sous-corpus de développement du corpus UD\_FRENCH-GSD (Guillaume *et al.*, 2019), deuxième plus grand corpus de français disponible dans UD. Celles-ci permettent à mettre en perspective, par ablation, l’impact des différents modèles de représentation lexicales. Elles nous ont également permis d’optimiser les hyperparamètres de nos modèles dans ces différentes configurations et de la procédure d’apprentissage. Ces optimisations ont été faites empiriquement, le nombre et le coût de ces expériences rendant une recherche systématique peu envisageable. Les résultats donnés pour ces expériences le sont sur le jeu de données de développement et non sur celui de test, afin d’éviter des effets de surapprentissage d’architecture.

Finalement nous testons sur les autres principaux corpus UD pour le français (UD\_FRENCH-SEQUOIA (Candito & Seddah, 2012; Bonfante *et al.*, 2018) et UD\_FRENCH-SPOKEN (Lacheret *et al.*, 2014; Gerdes & Kahane, 2017)) à l’exception du French Treebank, pour lequel nous reprenons la version utilisée pour la campagne d’évaluation SPMRL 2013 (Seddah *et al.*, 2013), ceci afin de pouvoir nous comparer plus facilement aux travaux existants. Tous les corpus UD sont utilisés dans leur version 2.7.

Nous avons utilisé principalement le modèle FLAUBERT (Le *et al.*, 2020) comme implémentation du modèle BERT. Pour la tâche d’analyse syntaxique nous avons remarqué lors d’expériences préliminaires qu’il se comporte généralement un peu mieux que le modèle CAMEMBERT (Martin *et al.*, 2020) même si les différences sont dans la marge d’erreur. Les modèles BERT utilisés sont les versions base-cased de FLAUBERT et de mBERT tels que distribués dans la version 4.2.2 de la bibliothèque TRANSFORMERS (Wolf *et al.*, 2020).

#### 3.1 Expériences de développement

Le modèle d’analyse syntaxique décrit ici se caractérise essentiellement par un enrichissement des représentations lexicales acquises à partir de gros volumes de données. L’ajout de ces représentations a également un coût non négligeable sur la taille du modèle et sur les temps d’exécution. Nous

TABLE 1 – Résultats (dev) des expériences d’ablations sur UD\_FRENCH-GSD 2.7

FASTTEXT	FLAUBERT	Lookup	Caractères	UPOS	UAS	LAS	CLAS
+	+	+	+	98,61	96,74	95,51	92,84
-	+	-	-	98,56	96,73	95,56	92,87
+	-	+	+	97,72	93,67	91,65	87,02
+	-	+	-	97,81	93,70	91,64	86,96
+	-	-	+	97,36	92,66	90,34	85,18
+	-	-	-	97,31	92,94	90,58	85,61
-	-	+	+	96,72	92,61	90,27	84,81
-	-	+	-	96,74	92,75	90,26	84,86
-	-	-	+	95,54	91,26	88,32	82,23

Lookup et caractères sont neutralisés en leur affectant des plongements de taille 2 (ce qui les rend moralement inexploitable par le modèle sans nécessiter un changement d’architecture), les plongements FLAUBERT sont neutralisés en les supprimant complètement des entrées. Enfin, pour les plongements FASTTEXT, – signifie que les plongements utilisés n’ont pas été préentraînés.

Les scores rapportés ici sont obtenus en répétant la procédure d’entraînement avec cinq germes aléatoires et en conservant les résultats du modèle le plus performant (en terme de LAS) sur le corpus de développement. Ces scores sont calculés par le script d’évaluation officiel de la campagne d’évaluation CoNLL 2018<sup>2</sup> dont nous conservons<sup>3</sup> les métriques UPOS, UAS, LAS, et CLAS (Nivre & Fang, 2017).

reportons ici quelques chiffres sur les performances de notre architecture en suivant une méthode d’ablation sur le corpus de développement. Les deux premières lignes du tableau 1 apportent un élément de réponse à la question : peut on se contenter de BERT comme seule représentation lexicale pour le parsing ? Il semble que la réponse est affirmative vu que la différence avec un modèle où toutes les autres représentations sont neutralisées est négligeable.

Notre seconde question est : peut-on se passer de BERT et obtenir des résultats équivalents en analyse syntaxique ? La troisième ligne du tableau donne notre meilleur modèle qui n’utilise pas de représentations de type BERT et contribue à apporter une réponse négative à la question.

En revanche, ce modèle utilise FASTTEXT et on remarque à partir des lignes 3 à 6 de la table que l’ablation additionnelle des représentations LOOKUP semble avoir un impact nettement plus significatif que l’omission des embeddings de caractères. Remarquons finalement à partir des lignes 7, 8 et 9 qu’utiliser des représentations FASTTEXT non-préappries entraîne une dégradation significative des performances. La suppression du sous modèle de caractères a là encore un effet négligeable.

Dans l’ensemble on remarque que l’apport de BERT est décisif : près de 5 points en LAS par rapport à un modèle appris uniquement à l’aide de représentations de mots sur le corpus d’entraînement. L’apport de FastText est plus modeste, et le sous-modèle de caractère semble avoir une contribution nulle dans les différentes configurations testées.

Ces tendances sont par ailleurs cohérentes pour les métriques UAS, LAS, mais aussi et surtout pour CLAS, ce qui suggère qu’elles sont bien liées à de réelles différences de performances pour la reconnaissance des structures syntaxiques à l’échelle de la phrase plutôt qu’à des différences dans

2. <https://github.com/ufal/conll2018/tree/e865d82e4c296d1660c9e7efcceb79aa418b6348>

3. Les métriques MLAS et BLEX, bien que plus récentes et préférées à CLAS dans cette campagne d’évaluation, n’auraient que peu de sens ici, puisque notre système ne prédit pas les traits morphologiques qu’elles mesurent.

TABLE 2 – Comparaisons entre les performances de nos modèles (sur différents corpus de test) en utilisant FLAUBERT, mBERT ou pas de représentations contextuelles (no BERT) et les performances de modèles à l’état de l’art.

(a) Résultats pour UD_FRENCH-GSD					(b) Résultats pour FTB-SPRML			
Modèle	UPOS	UAS	LAS	CLAS	Modèle	UPOS	UAS	LAS
mBERT	<i>98,13</i>	93,95	92,08	88,00	mBERT	<i>98,59</i>	91,68	88,48
FLAUBERT	<b>98,56</b>	95,67	<b>94,19</b>	<b>91,16</b>	FLAUBERT	<b>98,78</b>	<b>92,56</b>	<b>89,64</b>
no BERT	97,24	91,74	89,04	84,02	no BERT	98,03	88,54	84,54
<i>Martin et al.</i>	<i>98,18</i>	—	<i>92,57</i>	—	<i>Le et al.</i>	—	91,61	88,47
Stanza	97,30	91,38	89,05	84,38	<i>Constant et al.</i>	—	89,19	85,86
UD-Pipe 2	97,98	92,55	90,31	—				

  

(c) Résultats pour UD_FRENCH-SEQUOIA					(d) Résultats pour UD_FRENCH-SPOKEN				
Modèle	UPOS	UAS	LAS	CLAS	Modèle	UPOS	UAS	LAS	CLAS
mBERT	<i>99,01</i>	94,26	92,75	90,02	mBERT	<i>97,19</i>	83,93	78,13	71,05
FLAUBERT	<b>99,36</b>	95,68	<b>94,40</b>	92,12	FLAUBERT	<i>96,75</i>	86,00	80,46	73,97
no BERT	97,14	87,84	84,73	79,43	no BERT	92,62	77,87	69,78	61,03
<i>Martin et al.</i>	<i>99,29</i>	—	<i>94,20</i>	—	<i>Martin et al.</i>	<i>97,09</i>	—	<b>81,81</b>	—
Stanza	98,19	90,47	88,34	81,77	Stanza	95,49	75,82	70,71	62,13
UD-Pipe 2	99,32	94,88	93,81	—	UD-Pipe 2	<b>97,23</b>	86,27	<i>81,40</i>	—

Les scores pour nos modèles sont obtenus en répétant la procédure d’entraînement avec trois germes aléatoires et en conservant les résultats du modèle le plus performant (en terme de LAS) sur le corpus de développement. Ces scores sont calculés par le script d’évaluation officiel de la campagne d’évaluation CoNLL 2018<sup>4</sup> en considérant les métriques UPOS, UAS et LAS, ainsi que CLAS (Nivre & Fang, 2017) pour les corpus UD. Pour les métriques disponibles dans tout l’état de l’art, nous notons le meilleur résultat **en gras** et les résultats proches (moins de 0,5 points de différence) en *italiques*.

l’apprentissage des structures locales liées aux mots fonctionnels — puisque cette dernière métrique ne tiens précisément pas compte de ces dépendances.

### 3.2 Résultats de test

Nous présentons dans la table 2, une comparaison de nos résultats sur les jeux de données de test des principaux corpus de français à ceux de l’état de l’art — quand ils sont disponibles.

L’ensemble des modèles auxquels nous nous comparons sont également des variantes de l’algorithme d’analyse à arbre maximal couvrant introduite par Dozat & Manning (2017). Ils diffèrent essentiellement par les représentations lexicales utilisés. Notons toutefois que l’analyseur Stanza (Qi et al., 2020) réalise sa propre segmentation en mots alors que l’analyseur que nous présentons utilise la segmentation proposée par le treebank. Ses résultats ne sont donc pas parfaitement comparables avec les nôtres et sont donnés à titre indicatif. Tous les modèles utilisent en entrée des embeddings produits par une variante de BERT à l’exception de l’analyseur de Constant et al. (2013) qui n’est pas un analyseur à apprentissage profond et de Stanza. UDPipe 2 (Straka et al., 2019) utilise des

représentations BERT issues du modèle mBERT (Devlin *et al.*, 2019), mais contrairement à nous ces auteurs ne les ajustent pas durant l’entraînement de leur analyseur.

Les résultats obtenus par nos modèles, aussi bien pour l’étiquetage morphosyntaxique que pour l’analyse en dépendance sont généralement meilleurs que l’état de l’art, à l’exception des résultats pour UD\_FRENCH-SPOKEN. Pour ce dernier corpus, nous conjecturons que les résultats moins bons que l’état de l’art que nous obtenons s’expliquent au moins en partie par la petite taille et la plus grande irrégularité du corpus d’apprentissage, qui joue en la défaveur de nos modèles : ceux-ci présentant un plus grand nombre de paramètres que les modèles auxquels nous les comparons, leur tendance au surapprentissage est d’autant plus grande. Nous recommandons donc, pour des réutilisations de notre systèmes pour lesquelles les performances sur ce type de données seraient cruciales, une optimisation spécifique des hyperparamètres pour ce corpus afin de compenser cette baisse de performances.

Le tableau 3 donne par ailleurs des estimations des performances de ces modèles à l’inférence en terme de nombre de phrases traitées par seconde et d’occupation mémoire. On constate sans surprise que le passage sur GPU améliore grandement la vitesse, en particulier pour un modèle utilisant FlauBERT, mais que tous ces modèles restent utilisable pour des volumes de données moyens, même sur un ordinateur personnel.

TABLE 3 – Performances de nos modèles à l’inférence. Les valeurs données sont évaluées sur le corpus de test de UD\_French-GSD. Les valeurs « CPU » sont mesurées sur un ordinateur portable (processeur Intel Core i7-6500U, 2.50GHz en utilisant deux cœurs) et les valeurs « GPU » sur une machine dédiée au calcul (GPU NVidia GeForce RTX 3090). Les incertitudes sont calculées sur 10 répétitions

Modèle	Vitesse CPU (phrases/s)	Vitesse GPU (phrases/s)	Mémoire vive (GiB)
FLAUBERT	4,09 ± 0,06	20,32 ± 0,08	5,5
no BERT	22,96 ± 0,12	33,03 ± 0,79	3,5

## 4 Conclusion

Au final, on observe que les résultats de notre analyseur sont à l’état de l’art pour le français écrit tant l’analyse syntaxique que pour l’étiquetage morphosyntaxique. On obtient ce résultat en adaptant les embeddings contextuels de BERT appris sur un modèle spécifique au français (FLAUBERT). L’apport d’autres représentations comme FastText ou les plongements de caractères semble mineur. Pour le français parlé, il reste manifestement une marge de progression qui mérite un approfondissement.

Parmi les perspectives, nous envisageons de tester cette architecture d’analyse syntaxique sur des langues moins bien dotée, comme l’ancien français, ce qui est un contexte de travail où l’utilisation de modèles de plongements contextuels est beaucoup plus complexe à mettre en oeuvre, le manque de données linguistiques disponible rendant complexe le développement de représentations lexicales telles que celles utilisées ici.

## Remerciements

Ce travail bénéficié du soutien du projet ANR Profiterole (PROcessing Old French Instrumented TEXTs for the Representation Of Language Evolution), projet ANR-16-CE38-0010.



## Références

- BOJANOWSKI P., GRAVE E., JOULIN A. & MIKOLOV T. (2017). Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics*, **5**, 135–146. DOI : [10.1162/tacl\\_a\\_00051](https://doi.org/10.1162/tacl_a_00051).
- BONFANTE G., GUILLAUME B. & PERRIER G. (2018). *Application of Graph Rewriting to Natural Language Processing*, volume 1. ISTE Wiley.
- CANDITO M. & SEDDAH D. (2012). Le corpus Sequoia : annotation syntaxique et exploitation pour l'adaptation d'analyseur par pont lexical. In *Actes de la conférence conjointe JEP-TALN-RECITAL 2012*, volume 2 : Association pour le Traitement Automatique des Langues.
- CHU Y.-J. & LIU T.-H. (1965). On the shortest arborescence of a directed graph. *Scientia Sinica*, **14**, 1396–1400.
- CONSTANT M., CANDITO M. & SEDDAH D. (2013). The ligm-alpage architecture for the spmrl 2013 shared task : Multiword expression analysis and dependency parsing. In *Proceedings of the EMNLP Workshop on Statistical Parsing of Morphologically Rich Languages (SPMRL 2013) : shared task track*.
- DEVLIN J., CHANG M.-W., LEE K. & 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)*, Minneapolis, Minnesota : Association for Computational Linguistics.
- DOZAT T. & MANNING C. D. (2017). Deep biaffine attention for neural dependency parsing. In *Proceedings of ICLR*.
- EDMONDS J. (1967). Optimum branchings. *Journal of Research of the National Bureau of Standards*, **71B**(4), 233. DOI : [10.6028/jres.071b.032](https://doi.org/10.6028/jres.071b.032).
- GERDES K. & KAHANE S. (2017). Trois schémas d'annotation syntaxique en dépendance pour un même corpus de français oral : le cas de la macrosyntaxe. In *Actes de l'atelier sur les corpus annotés du français*.
- GUILLAUME B., DE MARNEFFE M.-C. & PERRIER G. (2019). Conversion et améliorations de corpus du français annotés en Universal Dependencies. *Traitement Automatique des Langues*, **60**(2), 71.
- HOCHREITER S. & SCHMIDHUBER J. (1997). Long short-term memory. *Neural computation*, **9**(8), 1735–1780.
- KINGMA D. P. & BA J. (2014). Adam : A Method for Stochastic Optimization. In *International Conference on Learning Representations*.
- LACHERET A., KAHANE S., BELIAO J., DISTER A., GERDES K., GOLDMAN J.-P., OBIN N., PIETRANDREA P. & TCHOBANOV A. (2014). Rhapsodie : a Prosodic-Syntactic Treebank for Spoken French. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation : European Language Resource Association*. HAL : [hal-00968959](https://hal.archives-ouvertes.fr/hal-00968959).
- LE H., VIAL L., FREJ J., SEGONNE V., COAVOUX M., LECOUTEUX B., ALLAUZEN A., CRABBÉ B., BESACIER L. & SCHWAB D. (2020). FlauBERT : Unsupervised language model pre-training for French. In *Proceedings of the 12th Language Resources and Evaluation Conference*, p. 2479–2490, Marseille, France : European Language Resources Association.



- MARTIN L., MULLER B., ORTIZ SUÁREZ P. J., DUPONT Y., ROMARY L., DE LA CLERGERIE É., SEDDAH D. & SAGOT B. (2020). CamemBERT : a tasty French language model. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, p. 7203–7219 : Association for Computational Linguistics.
- NIVRE J. & FANG C.-T. (2017). Universal Dependency Evaluation. In *Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017)*, p. 86–95 : Association for Computational Linguistics.
- QI P., ZHANG Y., ZHANG Y., BOLTON J. & MANNING C. D. (2020). Stanza : A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics : System Demonstrations*.
- SEDDAH D., TSARFATY R., KÜBLER S., CANDITO M., CHOI J. D., FARKAS R., FOSTER J., GOENAGA I., GOJENOLA GALLETEBEITIA K., GOLDBERG Y., GREEN S., HABASH N., KUHLMANN M., MAIER W., NIVRE J., PRZEPIÓRKOWSKI A., ROTH R., SEEKER W., VERSLEY Y., VINCZE V., WOLIŃSKI M., WRÓBLEWSKA A. & VILLEMONTÉ DE LA CLERGERIE E. (2013). Overview of the SPMRL 2013 Shared Task : A Cross-Framework Evaluation of Parsing Morphologically Rich Languages. In *Proceedings of the Fourth Workshop on Statistical Parsing of Morphologically-Rich Languages*, p. 146–182 : Association for Computational Linguistics.
- STRAKA M., STRAKOVÁ J. & HAJIC J. (2019). Evaluating contextualized embeddings on 54 languages in POS tagging, lemmatization and dependency parsing. *CoRR*, **abs/1908.07448**.
- VASWANI A., SHAZEER N., PARMAR N., USZKOREIT J., JONES L. & GOMEZ A. (2017). Attention is all you need. *Advances in neural information processing systems*, p. 5998–6008.
- WOLF T., DEBUT L., SANH V., CHAUMOND J., DELANGUE C., MOI A., CISTAC P., RAULT T., LOUF R., FUNTOWICZ M., DAVISON J., SHLEIFER S., VON PLATEN P., MA C., JERNITE Y., PLU J., XU C., LE SCAO T., GUGGER S., DRAME M., LHOEST Q. & RUSH A. (2020). Transformers : State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing : System Demonstrations*, p. 38–45 : Association for Computational Linguistics. DOI : [10.18653/v1/2020.emnlp-demos.6](https://doi.org/10.18653/v1/2020.emnlp-demos.6).
- ZEMAN D. *ET AL.* (2020). Universal dependencies 2.7. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.