



UNIVERSITÉ PARIS DESCARTES

PROJET PLURIDISCIPLINAIRE-MLSD

BIOMEDICAL NAMED ENTITY RECOGNITION

Team:

ABIDAR BOUCHRA

BOUSSEBAINE MUSTAPHA

LARBI ABDERRAHMANE

Teacher:

AFFELDT SEVERINE

2019-2020

CONTENTS

1	Introduction	3
2	Les algorithmes utilisés	4
2.1	CRF	4
2.1.1	CRF théoriquement	4
2.1.2	CRF techniquement	5
2.2	BiLSTM-CNN-CRF	5
2.2.1	BiLSTM-CNN-CRF théoriquement	5
2.2.2	BiLSTM-CNN-CRF techniquement	6
2.3	Bert	7
2.3.1	Flair	8
2.4	Spacy	10
3	Les données	11
4	Étude comparative	13
5	Modèle généralisé	15
6	Amélioration possibles	17
6.1	L'ajout des features	17
6.2	Changement d'architecture	17
6.3	Apprentissage continuée	17
7	conclusion	19
	References	20

INTRODUCTION

La reconnaissance d'entités nommées trouve leurs origines dans le domaine de la linguistique. Alors considérées comme sous-tâches de l'extraction d'information, ces disciplines ont rapidement attiré l'attention de différents domaines scientifiques tels que la biologie et la biomédecine. Avec le volume important des connaissances médicales numérisées à large échelle, retrouver automatiquement une information de haute précision est devenu un défi.

La reconnaissance biomédicale des entités nommées (Bio-NER) est une tâche fondamentale dans la gestion des termes textuels biomédicaux importants tels que les maladies, les gènes, les produits chimiques.

L'ensemble des études présente dans ce document, s'intéressent principalement à expérimenter les approches algorithmiques profonds en se basant sur les entités présentes sur différents corpus biomédicales.

Ce manuscrit est organisé selon quatre grandes parties. Nous commencerons par définir l'ensemble des algorithmes utilisés, ensuite les différents corpus expérimentés. Nous exploiterons ainsi les résultats obtenus sous forme d'une étude comparative. Puis nous présenterons un modèle que nous avons entraîné sur un corpus regroupant trois différentes thématiques. Ces dernières résumant l'ensemble des travaux réalisés jusqu'à présent. Néanmoins, nous finirons par proposer quelques techniques trouvées dans la littérature du NER pour améliorer la précision de la détection des entités nommées Biomédical.

LES ALGORITHMES UTILISÉS

2.1 CRF

2.1.1 CRF THÉORIQUEMENT

Les conditional Random Fields (ou CRF) sont des modèles Markoviens conditionnels qui ont été proposés pour remédier à certains défauts des modèles Markoviens plus classiques. Ils peuvent être appliqués à des données séquentielles, des données de type arbre et de façon plus générale à tout type de données structurées.

Les CRFs prennent comme entrées des données séquentielle avec une prise en compte de contexte précédent lors des prédictions sur un point de données. Nous modélisons mathématiquement cet aspect par "features de connaissances" définies comme suite :

$$f(X, i, l_{i-1}, l_i) \quad (1)$$

Le but de *features de connaissances* est de caractériser la séquence. Elle se base sur l'étiquette du mot précédent et du mot actuel. La construction du champ conditionnel est faite en attribuant à chaque *features de connaissances* un ensemble de poids (λ) que l'algorithme va apprendre. La fonction suivante explique la distribution de probabilité pour les champs aléatoires conditionnels CRF [1].

$$P(y, X, \lambda) = \frac{1}{Z(X)} \exp \sum_{i=1}^n \sum_j \lambda_j f_j(X, i, y_{i-1}, y_i) \quad (2)$$

et :

$$Z(X) = \sum_{y' \in y} \sum_{i=1}^n \sum_j \lambda_j f_j(X, i, y'_{i-1}, y'_i) \quad (3)$$

Les champs aléatoires conditionnels définissent les fonctions d'entité nécessaires, en initialisant les poids à des valeurs aléatoires, puis ils utilisent la descente de gradient de manière itérative jusqu'à la convergence des λ . L'équation de la mise à jour finale de descente de gradient pour

CRF est défini par l'équation (2) [1]:

$$\lambda = \lambda + \alpha \left[\sum_{k=1}^m F_j(y^k, x^k) + \sum_{k=1}^m P(y|x^k, \lambda) F_j(y, x^k) \right] \quad (4)$$

Les CRF sont similaires à la régression logistique, car ils utilisent la distribution de probabilité conditionnelle. La différence est que l'algorithme CRF utilise "features de connaissances" comme entrées séquentielles.

Pour résumer, nous présenterons une synthèse d'utilisation un CRF.

- Il faut fournir des exemples (x, y) .
- Il faut redéfinir sous forme de features ses connaissances.
- Il faut lancer l'apprentissage des K poids λ_k :
 1. En maximisant la log-vraisemblance
 2. Méthode : descente de gradient
- Une fois le modèle fixé, il permet de trouver l'annotation y qui maximise $p(y|x)$ pour tout nouveau x .

2.1.2 CRF TECHNIQUEMENT

Pour tester cette approche, nous avons expérimenté deux packages : *sklearn_crfsuite* et *pycrfsuit*. L'objectif est de pouvoir comparer les scores de performances des modèles.

Pour rajouter des informations aux différents corpus, nous avons utilisé le package *NLTK* pour avoir les entités POS. Les différentes expérimentations de CRF sur les corpus de thématiques différents (maladies, gènes et produits chimiques) existent dans le repo git: [code](#).

2.2 BiLSTM-CNN-CRF

2.2.1 BiLSTM-CNN-CRF THÉORIQUEMENT

Les études [2] ont montré que CNN est une approche efficace pour extraire des informations morphologiques (comme le préfixe ou le suffixe d'un mot) à partir de caractères de mots (figure 1).

Pour de nombreuses tâches d'étiquetage séquentiel, il est avantageux d'avoir accès aux deux contextes : passés et futurs. Cependant, une solution élégante dont l'efficacité a été prouvée par des travaux antérieurs [3] est le LSTM bidirectionnel (BLSTM).

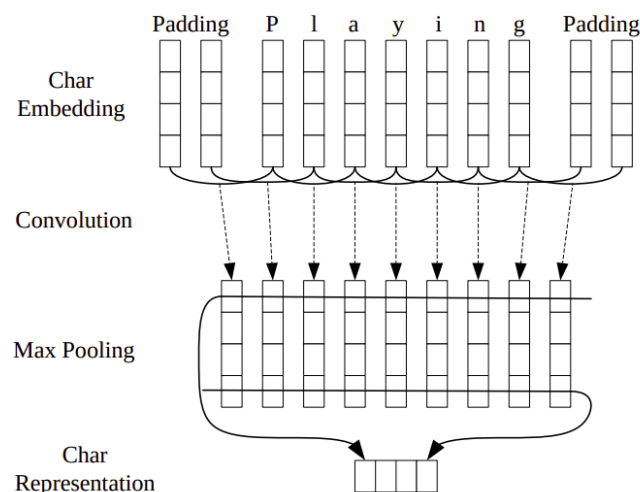


Figure 1: CNN pour l'extraction de représentations de mots au niveau des caractères.

L'idée de base est de présenter chaque séquence en avant et en arrière dans deux états cachés distincts afin de saisir les informations passées et futures, respectivement. Ensuite, les deux états cachés sont concaténés pour former le résultat final. La figure(2) illustre l'architecture d'un BiLSTM [4]. Vous pouvez consulter l'ensemble d'expérimentation réalisé par l'équipe dans le dépôt git [code](#).

Le modèle BiLSTM-CNN-CRF utilise ces deux architectures :

- Encodeur CNN pour la représentation embeddings.
- LSTM bidirectionnel pour le codage au niveau du mot.
- Champs aléatoires conditionnels (CRF) pour le décodage de sortie.

La figure(3) présente l'architecture de modèle BiLSTM-CNN-CRF.

2.2.2 BiLSTM-CNN-CRF TECHNIQUEMENT

Nous avons implémenté l'architecture BiLSTM CNN CRF en utilisant Keras Tensorflow. Nous avons utilisé le codage par CNN pour le codage par réseau neuronal à convolution pour la représentation des mots au niveau des caractères, le LSTM bidirectionnel pour le codage au niveau des mots et les champs aléatoires conditionnels (couche CRF) pour les décodages de sortie.

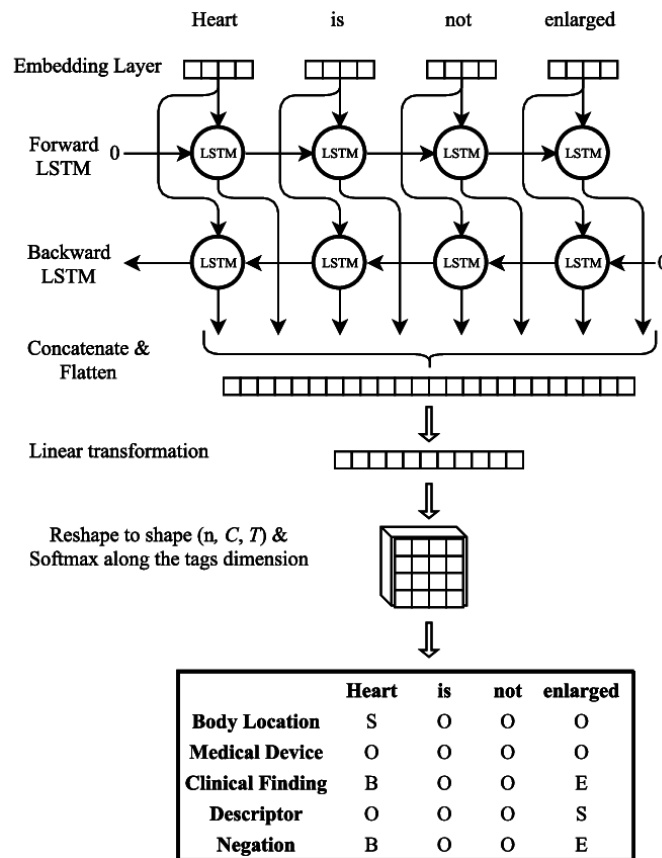


Figure 2: L'architecture de BLSTM

2.3 BERT

Bert a été présenté par Google. Ce modèle implémente plusieurs mécanismes notamment : *Transformer* et *Masked LM (MLM)*.

Transformer apprend les relations contextuelles entre les mots dans un texte[5]. Ce mécanisme contient deux parties : encodeur qui lit l'entrée de texte et décodeur pour prédire. Le mécanisme *Transformer* permet au modèle d'apprendre le contexte d'un mot en se basant sur l'ensemble de son environnement (gauche et droite du mot). Contrairement aux modèles directionnels, qui lisent l'entrée de texte séquentiellement (de gauche à droite ou de droite à gauche).

La deuxième stratégie qui rend Bert puissant est l'opération *Masked LM (MLM)*. Ce mécanisme permet de masquer 15% des mots de chaque séquence. Le modèle tente ensuite de prédire la valeur d'origine des mots masqués, sur la base du contexte (*Transformer*) fourni par les autres mots non masqués de la séquence (figure(4)).

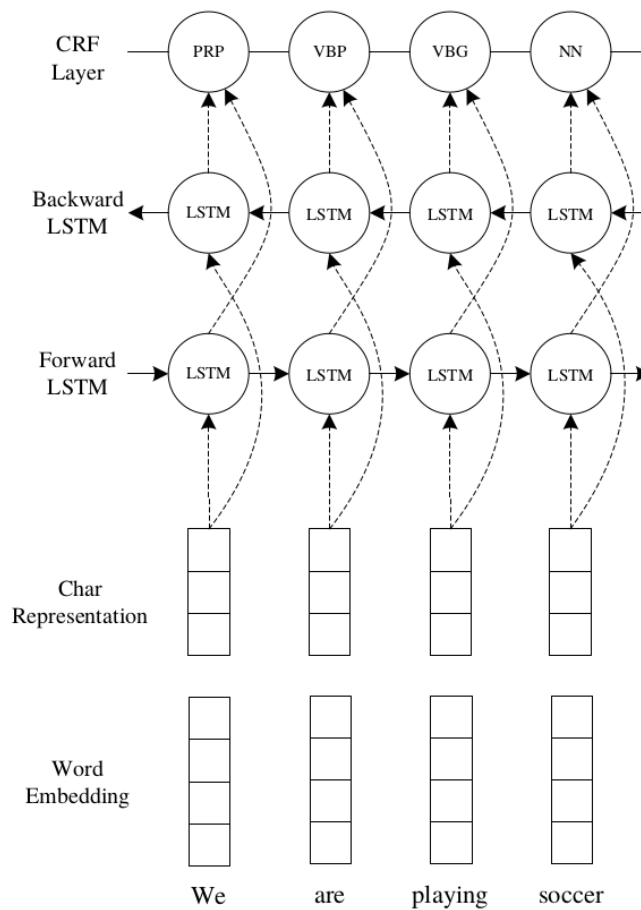


Figure 3: L'architecture de BLSTM-CNN-CRF[2]

2.3.1 FLAIR

Flair s'appuie sur une modélisation du langage neuronale au niveau des caractères pour apprendre des représentations contextuelles puissantes du langage humain à partir de grands corpus. La figure(5) présente l'architecture générale de Flair. Dans l'architecture de Flair, une phrase est vue en tant que séquence de caractères dans un modèle de langage de caractères bidirectionnel ML (Partie jaune de figure (5)) qui a été entraîné sur de très grands corpus de textes non étiquetés. À partir de ce LM, nous récupérons word a contextual embedding en enregistrant les états de la cellule du premier et du dernier caractère. Ce word embedding est ensuite passé dans un séquence BiLSTM-CRF vanille (bleu dans la figure 5).

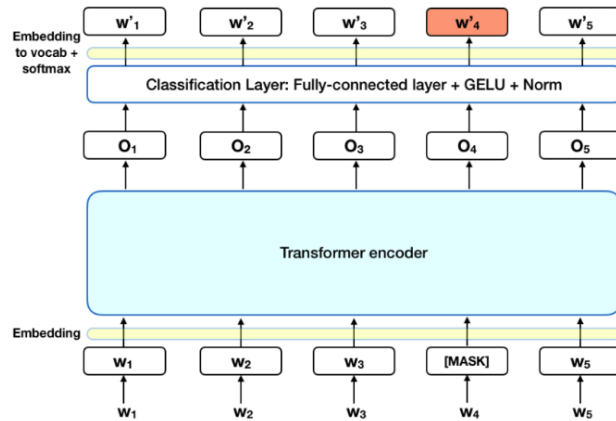


Figure 4: Mécanisme du Masked LM[5]

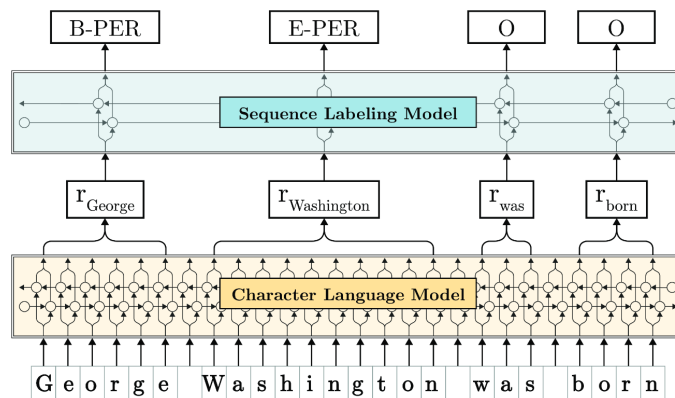


Figure 5: L'architecture de Flaire [6]

2.4 SPACY

Spacy est un outil pour la déction des NER. Il prend des décisions sur une base d'une prédiction entraînée pendant l'apprentissage. La figure suivant (figure 6) [10] explique le fonctionnement de spacy. Le format des données d'entraînement pour le modèle *Spacy* est différent de ceux

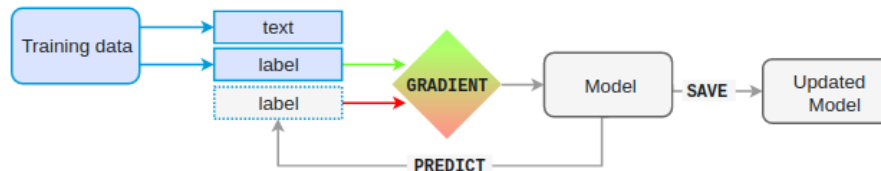


Figure 6: Le fonctionnement de spacy[10]

d'autres modèles. D'où la nécessité de procéder à une étape de pre-processing. Un fichier *helper.py* présent dans le repôt git [code](#) contient les différentes fonctions permettant de convertir automatiquement les données sous le format adapté aux inputs Spacy.

Dans l'utilisation de Spacy, la seule métrique qui évalue le modèle est *Loss*([resultat de loss](#)).

Les graphes suivants expliquent l'évolution de loss en fonction des epochs

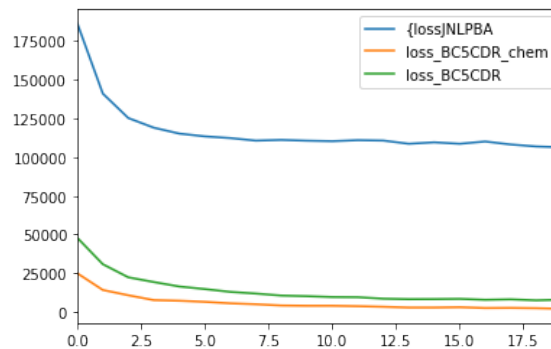


Figure 7: Loss spacy

LES DONNÉES

L'ensemble de données utilisées dans cette étude sont présentées dans le tableau ci-dessous (Table.1). Chaque corpus présente une thématique lié à un domaine bio-Médical notamment : les maladies, gènes, protéines et les produits chimiques.

DataSets					
corpus	Type d'Entité	Phrases	annotations	Taille	Publication
NCBI	Diseases	7639	6881	793	source
BC5CDR	Diseases	14228	12852	1500	source
JNLPBA	Genes	22562	35336	2404	source
BC5CDR-chem	Chemicals	14228	15935	1500	source

Table 1: Données d'entraînement

Les captures d'écrans suivantes présentent les top features dans les corpus (*NCBI*, *BC5CDR*, *BC5CDR-chemical* et *JNLPBA*) avec une estimation des poids associées (code)

yB-Disease top features		yE-Disease top features		yI-Disease top features		yS-Disease top features	
Weight	Feature	Weight	Feature	Weight	Feature	Weight	Feature
+5.483	-1 word lower/hypersensitivity	+7.082	-1 word lower/palae	+6.146	-1 word lower/hypersensitivity	+6.688	-1 word lower/hypoglycemia
+5.119	-1 word lower/dysphagia	+6.131	-1 word lower/dysphagia	+5.897	-1 word lower/epididym	+5.840	-1 word lower/tumors
+4.239	-1 word lower/dysplasia	+5.289	-1 word lower/cancers	+5.530	-1 word lower/rhinit	+5.600	-1 word lower/cataracts
+4.217	-1 word lower/prur	+4.842	-1 word lower/tetanus	+5.132	-1 word lower/gastric	+5.434	-1 word lower/asthma
+4.195	-1 word lower/lymphed	+4.810	-1 word lower/hearing	+5.047	-1 word lower/achilles	+5.291	-1 word lower/3 mia
+4.187	-1 word lower/hypopleas	+4.471	-1 word lower/syndrome	+4.866	-1 word lower/syndrome	+5.278	-1 word lower/neurodegeneration
+4.173	-1 word lower/disease	+4.414	-1 word lower/veins	+4.502	-1 word lower/fingers	+5.010	-1 word lower/hypomelanisation
+4.162	-1 word lower/contractures	+4.195	-1 word lower/wasting	+4.469	-1 word lower/breast	+4.926	-1 word lower/arrhythm
+4.039	-1 word lower/von	+4.142	-1 word lower/3 mia	+4.291	-1 word lower/linked	+4.803	-1 word lower/potter
+4.018	-1 word lower/deficiencies	+4.081	-1 word lower/pigments	+4.210	-1 word lower/syndrome	+4.825	-1 word lower/obesity
+3.967	-1 word lower/demeylation	+4.105	-1 word lower/anchored	+4.154	-1 word lower/von	+4.699	-1 word lower/oste
+3.926	-1 word lower/rehydration	+4.081	-1 word lower/epididym	+4.129	-1 word lower/component	+4.329	-1 word lower/3 gla
+3.858	-1 word lower/defects	+3.950	-1 word lower/von	+4.118	-1 word lower/methachol	+4.269	-1 word lower/tumour
+3.796	-1 word lower/cause	+3.901	-1 word lower/disorder	+4.110	-1 word lower/von	+4.211	-1 word lower/hypoglycemia
+3.791	-1 word lower/ewing	+3.892	-1 word lower/tolerance	+4.081	-1 word lower/oes	+4.118	-1 word lower/3 mia
+3.769	-1 word lower/status	+3.835	-1 word lower/neoplasms	+4.008	-1 word lower/palae	+4.106	-1 word lower/cataract
+3.769	-1 word lower/cerebral	+3.787	-1 word lower/wilebrand	+4.000	-1 word lower/ovarian	+4.083	-1 word lower/infertility
+3.769	-1 word lower/disorders	+3.772	-1 word lower/atriumvent	+3.938	-1 word lower/dysmyelinating	+4.075	-1 word lower/obese
+3.740	-1 word lower/epiphyseal	+3.765	-1 word lower/linked	+3.928	-1 word lower/inborn	+4.047	-1 word lower/3 AMN
+3.667	-1 word lower/deficiency	+3.749	-1 word lower/ovarian	+3.918	-1 word lower/anchored	+4.047	-1 word lower/arm
+430 more positive		+430 more positive		+430 more positive		+430 more positive	
141 more negative		141 more negative		148 more negative		126 more negative	

Figure 8: NCBI top features

yB-Chemical top features		yI-Chemical top features		yS-Chemical top features	
Weight	Feature	Weight	Feature	Weight	Feature
+7.756	-1 word lower/antidepressant	+7.422	-1 word lower/vitamin	+6.743	-1 word lower/vomiting
+7.359	-1 word lower/bupropion	+6.584	-1 word lower/cad	+6.505	-1 word lower/infection
+6.889	-1 word lower/3 cin	+6.935	-1 word lower/seizure	+4.984	-1 word lower/3 oma
+6.273	-1 word lower/jumal	+6.486	-1 word lower/3 methyl	+4.804	-1 word lower/3 mia
+6.134	-1 word lower/antidepressants	+4.478	-1 word lower/benies	+4.896	-1 word lower/emollic
+5.887	-1 word lower/coude	+4.476	-1 word lower/benies	+4.819	-1 word lower/disorder
+5.780	-1 word lower/3 sin	+4.423	-1 word lower/lophavir	+4.899	-1 word lower/urter
+5.668	-1 word lower/3 iol	+4.277	-1 word lower/contraceptives	+4.488	-1 word lower/disease
+5.667	-1 word lower/3 iol	+4.248	-1 word lower/3 iole	+4.463	-1 word lower/urter
+5.565	-1 word lower/contraceptives	+4.218	-1 word lower/dactinomycin	+4.474	-1 word lower/coaction
+5.532	-1 word lower/3 iol	+4.083	-1 word lower/toxicity	+4.315	-1 word lower/3 iol
+5.435	-1 word lower/3 pam	+4.080	-1 word lower/carboxylated	+4.276	-1 word lower/seizures
+5.422	-1 word lower/3 iol	+3.982	-1 word lower/pregnan	+4.199	-1 word lower/infect
+5.018	-1 word lower/3 iol	+3.853	-1 word lower/alkylating	+4.135	-1 word lower/midline
+4.762	-1 word lower/phenobarbital	+3.840	-1 word lower/3 iol	+4.105	-1 word lower/midline
+4.747	-1 word lower/glutamate	+3.799	-1 word lower/3 iol	+4.032	-1 word lower/syndrome
+4.703	-1 word lower/3 iol	+3.772	-1 word lower/3 iol	+3.996	-1 word lower/seizures
+4.660	-1 word lower/3 iol	+3.787	-1 word lower/3 iol	+3.937	-1 word lower/vomiting
+4.578	-1 word lower/bilirubin	+3.775	-1 word lower/3 iol	+3.972	-1 word lower/creal
+179 more positive		+179 more positive		+179 more positive	
179 more negative		110 more negative		227 more negative	

Figure 10: BC5CDR chemical top features

yB-DNA top features		yE-DNA top features		yI-DNA top features		yS-DNA top features	
Weight	Feature	Weight	Feature	Weight	Feature	Weight	Feature
+5.842	-1 word lower/genes	+6.054	-1 word lower/g22	+5.977	-1 word lower/genes	+5.989	-1 word lower/intons
+5.465	-1 word lower/many	+5.680	-1 word lower/2	+5.940	-1 word lower/sequences	+5.731	-1 word lower/transcripts
+5.080	-1 word lower/genes	+5.553	-1 word lower/pro-oncogenes	+5.264	-1 word lower/promoters	+5.392	-1 word lower/c-site
+4.953	-1 word lower/promoter-reg	+5.461	-1 word lower/pro	+5.193	-1 word lower/elements	+5.165	-1 word lower/protoc
+4.728	-1 word lower/genes	+5.362	-1 word lower/g14	+5.145	-1 word lower/histocompatibility	+4.808	-1 word lower/gemine
+4.715	-1 word lower/respective	+4.950	-1 word lower/g32	+5.027	-1 word lower/c-rich	+4.885	-1 word lower/gemine
+4.704	-1 word lower/genes	+4.840	-1 word lower/g14	+4.989	-1 word lower/histocompatibility	+4.812	-1 word lower/3 iol
+4.543	-1 word lower/recombination	+4.493	-1 word lower/purine-rich	+4.887	-1 word lower/site-like	+4.682	-1 word lower/chromatin
+4.513	-1 word lower/genes	+4.461	-1 word lower/genes	+4.641	-1 word lower/elements	+4.602	-1 word lower/chromatin
+4.505	-1 word lower/promoters	+4.430	-1 word lower/antipromoters	+4.806	-1 word lower/element	+4.441	-1 word lower/expression
+4.484	-1 word lower/cna	+4.378	-1 word lower/3 iol	+4.767	-1 word lower/nucleotides	+4.445	-1 word lower/direct
+4.362	-1 word lower/3 iol	+4.361	-1 word lower/element	+4.764	-1 word lower/induction	+4.379	-1 word lower/induction
+4.305	-1 word lower/plasmids	+4.320	-1 word lower/constructs	+4.703	-1 word lower/g22	+4.358	-1 word lower/rhombin
+4.291	-1 word lower/ind	+4.287	-1 word lower/genes	+4.599	-1 word lower/constructs	+4.382	-1 word lower/afect
+4.263	-1 word lower/promoter	+4.263	-1 word lower/p27kip1	+4.577	-1 word lower/g22	+4.338	-1 word lower/3 iol
+4.250	-1 word lower/elements	+4.238	-1 word lower/vectors	+4.551	-1 word lower/vectors	+4.260	-1 word lower/3 iol
+4.178	-1 word lower/sequences	+4.225	-1 word lower/3 iol	+4.522	-1 word lower/constructs	+4.201	-1 word lower/enhancers
+4.161	-1 word lower/early	+4.222	-1 word lower/constructs	+4.419	-1 word lower/kapap-mutated	+4.119	-1 word lower/3-region
+4.136	-1 word lower/3 iol	+4.200	-1 word lower/3 iol	+4.391	-1 word lower/3 iol	+4.045	-1 word lower/3 iol
+4.107	-1 word lower/prior	+4.243 more positive		+4.389 more positive		+4.045	-1 word lower/3 iol
3638 more positive		295 more positive		378 more positive		422 more positive	
273 more negative		378 more negative		378 more negative		115 more negative	

yB-cell_line top features		yE-cell_line top features		yI-cell_line top features		yS-cell_line top features	
Weight	Feature	Weight	Feature	Weight	Feature	Weight	Feature
+5.583	-1 word lower/lines	+5.420	-1 word lower/line	+6.132	-1 word lower/lines	+6.038	-1 word lower/thp-1
+5.307	-1 word lower/different	+5.191	-1 word lower/line	+6.081	-1 word lower/line	+5.403	-1 word lower/promyelocytes
+5.018	-1 word lower/clones	+5.120	-1 word lower/lines	+4.783	-1 word lower/clones	+4.725	-1 word lower/heterokaryons
+4.787	-1 word lower/colonies	+5.031	-1 word lower/colonies	+4.713	-1 word lower/line	+4.126	-1 word lower/huvecs
+4.712	-1 word lower/several	+4.810	-1 word lower/line	+4.650	-1 word lower/colonies	+4.002	-1 word lower/3 iol
+4.678	-1 word lower/primed	+4.414	-1 word lower/population	+4.635	-1 word lower/sk-b-e	+4.061	-1 word lower/1
+4.617	-1 word lower/various	+4.410	-1 word lower/lines	+4.572	-1 word lower/line	+3.994	-1 word lower/3 iol
+4.560	-1 word lower/other	+4.366	-1 word lower/b2	+4.385	-1 word lower/histocompatibility	+3.956	-1 word lower/2
+4.406	-1 word lower/same	+4.342	-1 word lower/h-5f	+4.318	-1 word lower/tcrpha	+3.863	-1 word lower/molt-4
+4.332	-1 word lower/transfectant	+4.171	-1 word lower/subclone	+4.276	-1 word lower/mousepromoted	+3.848	-1 word lower/j-1
+4.331	-1 word lower/line	+4.162	-1 word lower/hybridomas	+4.164	-1 word lower/lines	+3.842	-1 word lower/mat
+4.290	-1 word lower/cells	+4.125	-1 word lower/clone 12	+4.143	-1 word lower/transformed	+3.722	-1 word lower/lines
+4.062	-1 word lower/tumorigenic	+4.067	-1 word lower/bearing	+4.067	-1 word lower/hybridoma	+3.732	-1 word lower/1CLs
+4.020	-1 word lower/epo-independent	+4.005	-1 word lower/b-cells	+4.012	-1 word lower/phytohemagglutinin	+3.684	-1 word lower/3 iol
+3.988	-1 word lower/pickle	+3.806	-1 word lower/transfectants	+3.965	-1 word lower/hybrids	+3.616	-1 word lower/hybrid
+3.903	-1 word lower/producing	+3.768	-1 word lower/culture	+3.859	-1 word lower/3 iol	+3.612	-1 word lower/transfectants
+3.825	-1 word lower/distinct	+3.723	-1 word lower/h104	+3.846	-1 word lower/hes-2-specific	+3.553	-1 word lower/hela
+3.780	-1 word lower/h40	+3.723	-1 word lower/junkat	+3.820	-1 word lower/hybridoma	+3.540	-1 word lower/3 iol
+3.771	-1 word lower/c3d4-thy-1	+3.685	-1 word lower/line	+3.819	-1 word lower/overexpressing	+3.523	-1 word lower/3 iol
+3.751	-1 word lower/preactivated	+3.685	-1 word lower/line	+3.819	-1 word lower/overexpressing	+3.523	-1 word lower/3 iol
1957 more positive		195 more positive		262 more positive		548 more positive	
145 more negative		195 more negative		38 more negative		38 more negative	

yB-protein top features		yE-protein top features		yI-protein top features		yS-protein top features	
Weight	Feature	Weight	Feature	Weight	Feature	Weight	Feature
+6.177	-1 word lower/rb-y	+3.344	-1 word lower/c-fos-positive	+6.307	-1 word lower/superfamily	+6.023	-1 word lower/ident
+6.487	-1 word lower/other	+3.328	-1 word lower/c-fos-positive	+5.462	-1 word lower/product	+5.809	-1 word lower/cytokines
+5.426	-1 word lower/c-fos-positive	+3.288	-1 word lower/sdk	+5.398	-1 word lower/transducers	+5.306	-1 word lower/alpha-helices
+5.283	-1 word lower/distinct	+3.020	-1 word lower/bsagp	+5.377	-1 word lower/receptors	+4.896	-1 word lower/lymphokines
+5.045	-1 word lower/factors	+3.823	-1 word lower/k41	+4.785	-1 word lower/of	+4.879	-1 word lower/lymphocytes
+4.969	-1 word lower/receptors	+3.439	-1 word lower/waf1/cip1	+4.728	-1 word lower/chains	+4.838	-1 word lower/24-hydroxylase
+4.911	-1 word lower/express	+3.427	-1 word lower/sk	+4.642	-1 word lower/complexes	+4.621	-1 word lower/enzyme
+4.749	-1 word lower/sphingomyelinase	+4.826	-1 word lower/docking	+4.604	-1 word lower/products	+4.621	-1 word lower/interferons
+4.718	-1 word lower/synthesis	+4.731	-1 word lower/activators	+4.534	-1 word lower/activators	+4.442	-1 word lower/3 iol
+4.703	-1 word lower/kineses	+4.601	-1 word lower/secreted	+4.511	-1 word lower/aminoc	+4.424	-1 word lower/p2c2-bis
+4.637	-1 word lower/m2	+4.503	-1 word lower/transporter	+4.485	-1 word lower/anti-cd70	+4.293	-1 word lower/3-region
+4.628	-1 word lower/complex	+4.467	-1 word lower/cytokines	+4.461	-1 word lower/family	+4.309	-1 word lower/3 iol
+4.573	-1 word lower/mab	+4.426	-1 word lower/c	+4.468	-1 word lower/pep2b/eamh11	+4.297	-1 word lower/pep2b/eamh11
+4.548	-1 word lower/different	+4.411	-1 word lower/122c	+4.464	-1 word lower/complexes	+4.265	-1 word lower/hybridoma
+4.547	-1 word lower/bnd	+4.398	-1 word lower/2i6-myc	+4.436	-1 word lower/antihemorrhage	+4.223	-1 word lower/3 iol
+4.514	-1 word lower/superfamily	+4.392	-1 word lower/122i5	+4.349	-1 word lower/125-125	+4.198	-1 word lower/signal-transduction
+4.448	-1 word lower/described	+4.387	-1 word lower/risk	+4.296	-1 word lower/transducers	+4.196	-1 word lower/zap-70
+4.402	-1 word lower/tetrapeptide	+4.357	-1 word lower/kinase-2	+4.284	-1 word lower/chimeras	+4.167	-1 word lower/cytokine
+4.380	-1 word lower/predominant	+4.271	-1 word lower/p38c-jun	+4.249	-1 word lower/stem-binding	+4.157	-1 word lower/pathway
+4.312	-1 word lower/acid-sensitive	+4.228	-1 word lower/substrates	+4.228	-1 word lower/3 iol	+4.108	-1 word lower/3 iol
482 more positive		417 more positive		415 more positive		479 more positive	
48 more negative		394 more positive		395 more positive		395 more positive	

Figure 11: Genes top features

yB-Chemical top features		yI-Chemical top features		yB-Disease top features		yI-Disease top features	
Weight	Feature	Weight	Feature	Weight	Feature	Weight	Feature
+7.756	-1 word lower/antidepressant	+7.422	-1 word lower/vitamin	+6.743	-1 word lower/migraine	+6.130	-1 word lower/vomiting
+7.359	-1 word lower/bupropion	+6.584	-1 word lower/cad	+6.505	-1 word lower/infection	+5.840	-1 word lower/tumors
+6.889	-1 word lower/3 cin	+6.935	-1 word lower/seizure	+4.984	-1 word lower/3 oma	+5.434	-1 word lower/asthma
+6.273	-1 word lower/jumal	+6.486	-1 word lower/3 methyl	+4.804	-1 word lower/3 mia	+5.291	-1 word lower/3 mia
+6.134	-1 word lower/antidepressants	+4.478	-1 word lower/benies	+4.896	-1 word lower/emollic	+4.819	-1 word lower/disorder
+5.887	-1 word lower/coude	+4.476	-1 word lower/benies	+4.819	-1 word lower/disorder	+4.699	-1 word lower/urter
+5.780	-1 word lower/3 sin	+4.423	-1 word lower/lophavir	+4.899	-1 word lower/urter	+4.483	-1 word lower/urter
+5.668	-1 word lower/3 iol	+4.277	-1 word lower/contraceptives	+4.488	-1 word lower/hypertensive	+4.483	-1 word lower/urter
+5.667	-1 word lower/3 iol	+4.248	-1 word lower/3 iole	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+5.565	-1 word lower/contraceptives	+4.218	-1 word lower/dactinomycin	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+5.532	-1 word lower/3 iol	+4.083	-1 word lower/toxicity	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+5.435	-1 word lower/3 pam	+4.080	-1 word lower/carboxylated	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+5.422	-1 word lower/3 iol	+3.982	-1 word lower/pregnan	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+5.018	-1 word lower/3 iol	+3.853	-1 word lower/alkylating	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+4.762	-1 word lower/phenobarbital	+3.840	-1 word lower/3 iol	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+4.747	-1 word lower/glutamate	+3.799	-1 word lower/3 iol	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+4.703	-1 word lower/3 iol	+3.772	-1 word lower/3 iol	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+4.660	-1 word lower/3 iol	+3.787	-1 word lower/3 iol	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+4.578	-1 word lower/bilirubin	+3.775	-1 word lower/3 iol	+4.474	-1 word lower/coaction	+4.474	-1 word lower/coaction
+179 more positive		+179 more positive		+179 more positive		+179 more positive	
179 more negative		110 more negative		227 more negative		126 more negative	

ÉTUDE COMPARATIVE

La table(2) résume les performances de détection d'entités de chaque algorithme. En se basant sur la dispersion des classes, nous avons utilisé le *F1 Score* comme un indice de performance. Dans Les 4 corpus testés, FLair renvoie le meilleur score.

Table 2: La performances des modèles

Algorithmes		F1 score(%)
CRF	NCBI	77
	BC5CDR	75
	JNLPBA	70
	BC5CDR-chem	76
BiLSTM-CRF	NCBI	70
	BC5CDR	65
	JNLPBA	
	BC5CDR-chem	60
Flair	NCBI	88
	BC5CDR	88
	JNLPBA	81
	BC5CDR-chem	93

L'algorithme *BiLSTM-CRF* a donné les faibles scores pour la détection des produits chimiques et les maladies sur le corpus BC5CDR. Nous proposons dans la section 6 "**Amélioration possibles**" des techniques trouvées dans la littérature pour augmenter la performance. Suite à un problème au niveau des ressources, nous n'avons pas réussi à entraîner le modèle BiLSTM-CRF sur le corpus *JNLPBA*.

Une tentative a été testée en diminuant la taille des corpus. Cette tentative n'a pas abouti aux résultats attendus : F1 score égale à 25%

L'exemple d'exécution des modèles sur des phrases différentes en temps réel est présenté dans la figure (12)

```
In [52]: print(doc_JNLPBA)
print("-----")
print("Tagged : ",NER_JNLPBA.predict(doc_JNLPBA)[0].to_tagged_string())

gene expression and NF-kappa B activation through CD28 requires reactive oxygen production by 5-
lipoxigenase .
-----
Tagged : gene <E-DNA> expression and NF-kappa <B-protein> B <E-protein> activation through CD2
8 <S-protein> requires reactive oxygen production by 5-lipoxygenase <S-protein> .

In [53]: print("Sentence :",doc_BC5CDR_chem)
print("-----")
print("Tagged : ",NER_BC5CDR_chem.predict(doc_BC5CDR_chem)[0].to_tagged_string())

Sentence : OBJECTIVES : The United Kingdom Parkinson ' s Disease Research Group ( UKPDGRG ) tria
l found an increased mortality in patients with Parkinson ' s disease ( PD ) randomized to recei
ve 10 mg selegiline per day and L - dopa compared with those taking L - dopa alone .
-----
Tagged : OBJECTIVES : The United Kingdom Parkinson ' s Disease Research Group ( UKPDGRG ) trial
found an increased mortality in patients with Parkinson ' s disease ( PD ) randomized to receive
10 mg selegiline <S-Chemical> per day and L <B-Chemical> - <I-Chemical> dopa <E-Chemical> compar
ed with those taking L <B-Chemical> - <I-Chemical> dopa <E-Chemical> alone .

In [54]: print(doc_BC5CDR)
print("-----")
print("Tagged : ",NER_BC5CDR.predict(doc_BC5CDR)[0].to_tagged_string())

In the in vivo study , the administration ( 50 mg / kg , i . p . ) of TET and FAN in mice showed
the inhibition of thrombosis by 55 % and 35 % , respectively , while acetylsalicylic acid ( ASA
, 50 mg / kg , i . p . ) , a positive control , showed only 30 % inhibition .
-----
Tagged : In the in vivo study , the administration ( 50 mg / kg , i . p . ) of TET <B-Chemical
> and FAN <B-Chemical> in mice showed the inhibition of thrombosis <B-Disease> by 55 % and 35 %
, respectively , while acetylsalicylic <B-Chemical> acid <I-Chemical> ( ASA <B-Chemical> , 50 mg
/ kg , i . p . ) , a positive control , showed only 30 % inhibition .
```

Figure 12: tester les modèles sur des phrases différentes

MODÈLE GÉNÉRALISÉ

Dans le but de pousser notre étude et expérimenter sur un corpus contenant les différentes thématiques abordées auparavant, nous avons décidé de créer un modèle entraîné sur un corpus regroupant les quatre thématiques. Le but est d'avoir un unique modèle capable de détecter les différents tags notamment : maladies, gènes et produits chimiques.

Le modèle Flair a été retenu pour cette tâche, suite à son résultat pertinent dans les corpus précédents ainsi que l'intuitive de l'architecture de ce dernier.

EPOCH	LOSS	PRECISION	RECALL	F1
1	6.93	0.6386	0.4443	0.5240
10	3.22	0.7955	0.6725	0.7289
20	2.80	0.7664	0.7846	0.7754
50	2.54	0.7812	0.7929	0.7870
72	2.53	0.7820	0.7955	0.7887

Table 3: Résultats d'entraînement du modèle générale

Avec l'entraînement du modèle sur ce corpus généralisé, nous avons atteint un score F1 de 79%. Un score qui reste intéressant en prenant en considération la distribution des classes non équilibrées ainsi que le nombre important de tagges différents présents dans le corpus construit.

Une autre manière de tester le modèle généralisé est d'essayer de tagger les features dans les mêmes phrases de figure(12). Ce que nous proposons dans la figure (13) `code`.

Sentence : gene expression and NF-kappa B activation through CD28 requires reactive oxygen production by 5-lipoxygenase .

Tagged : gene expression and NF-kappa <B-protein> B <E-protein> activation through CD28 <S-protein> requires reactive oxygen production by 5-lipoxygenase .

Sentence : In the in vivo study , the administration (50 mg / kg , i . p .) of TET and FAN in mice showed the inhibition of thrombosis by 55 % and 35 % , respectively , while acetylsalicylic acid (ASA , 50 mg / kg , i . p .) , a positive control , showed only 30 % inhibition .

Tagged : In the in vivo study , the administration (50 mg / kg , i . p .) of TET <B-Chemical> and FAN <B-Chemical> in mice showed the inhibition of thrombosis <B-Disease> by 55 % and 35 % , respectively , while acetylsalicylic <B-Chemical> acid <I-Chemical> (ASA <B-Chemical> , 50 mg / kg , i . p .) , a positive control , showed only 30 % inhibition .

Sentence : OBJECTIVES : The United Kingdom Parkinson ' s Disease Research Group (UKPDRG) trial found an increased mortality in patients with Parkinson ' s disease (PD) randomized to receive 10 mg selegiline per day and L - dopa compared with those taking L - dopa alone .

Tagged : OBJECTIVES : The United Kingdom Parkinson <B-Disease> ' <I-Disease> s <I-Disease> Disease <I-Disease> Research Group (UKPDRG) trial found an increased mortality in patients with Parkinson <B-Disease> ' <I-Disease> s <I-Disease> disease <I-Disease> (PD <B-Disease>) randomized to receive 10 mg selegiline <B-Chemical> per day and L <B-Chemical> - <I-Chemical> dopa <I-Chemical> compared with those taking L <B-Chemical> - <I-Chemical> dopa <I-Chemical> alone .

Figure 13: Exemples des phrases tagger avec le modèle généralisé

AMÉLIORATION POSSIBLES

6.1 L'AJOUT DES FEATURES

Nous pouvons toujours trouver dans la littérature des techniques pour améliorer le score de la performance des modèles. L'ensemble des travaux effectués dans l'article [7], montre une réussite remarquable d'incrémenter le score F1 en ajoutant des features pour améliorer les performances d'un CRF. Ce que nous pouvons faire aussi dans notre cas.

Une autre approche consiste à faire varier l'architecture de BiLSTM-CRF pour essayer d'augmenter le degré de précision de moddèle. Ce dépôt git présente un ensemble des architectures [code](#)

6.2 CHANGEMENT D'ARCHITECTURE

L'une des architectures qui a donné des scores de performance remarquable environ 90% est BERT-BiLSTM-CRF-NER. Le lien suivant illustre les différents travaux présentés par les chercheurs [code](#). La figure (14) présente l'architecture de ce modèle.

6.3 APPRENTISSAGE CONTINUÉE

Une autre approche proposée pour augmenter la performance de la détection des NER consiste à créer un **Meta-Modèle** qui regroupe le modèle *Flair*, *Spacy*, *CRF* et *BLSTM*. L'idée derrière ce regroupement est de donner la possibilité à chaque modèle d'apprendre en permanence en fonction des résultats des autres modèles. Nous pouvons mettre en place cette approche en implémentant les étapes suivantes :

- Le modèle prend comme entrée une phrase à tagger.
- La phrase en question passe par les trois modèles: *Flair*, *Spacy*, *CRF* et *BLSTM*
- Chaque modèle renvoie ces outputs.

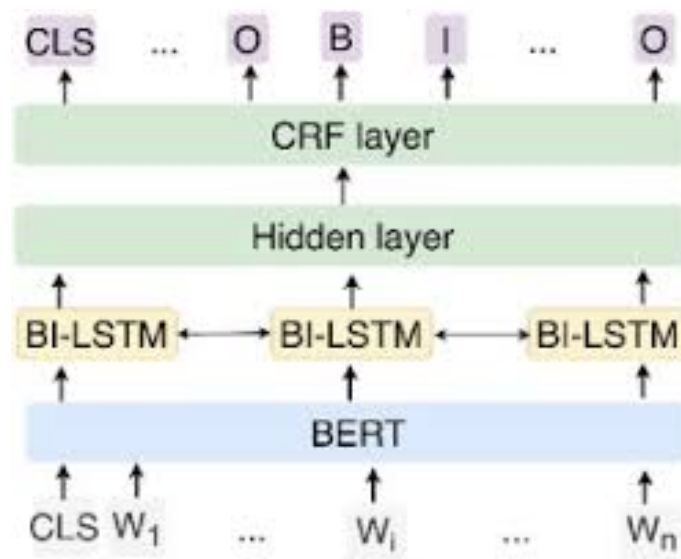


Figure 14: L'architecture Bert-CRF-BLSTM [9]

- L'algorithme sélectionne le modèle qui renvoie le meilleur tag.
- L'algorithme lance un reentrainement en utilisant l'output génère dans l'étape précédent comme donnée d'entraînement pour les autres modèles.

CONCLUSION

Nous avons eu l'occasion de traiter une problématique très répandue dans les domaines de recherche de data science, plus précisément le domaine de NLP dans un cadre BioMedical. Durant ce projet, nous étions amenées à réaliser un NER sur des différentes thématiques BioMedical.

La richesse de ce domaine, nous a permis de découvrir, utiliser et expérimenter plusieurs architectures d'apprentissage automatique profond.

Nous avons constaté qu'il est toujours possible de combiner plusieurs architectures adaptées pour améliorer les performances des modèles. Une tâche importante d'un data scientist

Nous avons constaté qu'il est toujours possible de combiner plusieurs architectures adaptées pour améliorer les performances des modèles.

REFERENCES

- [1] John Lafferty, Andrew McCallum, Fernando Pereira *Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data*. In Proceedings of the Eighteenth International Conference on Machine Learning (ICML-2001), 2001
- [2] Xuezhe Ma and Eduard Hovy *End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF*, 2016
- [3] Eliyahu Kiperwasser, Yoav Goldberg *Simple and Accurate Dependency Parsing Using Bidirectional LSTM Feature Representations*, 2016
- [4] Savelie Cornegruta, Robert Bakewell, Samuel Withey and Giovanni Montana, *Modelling Radiological Language with Bidirectional Long Short-Term Memory Networks*, 2016
- [5] BERT Explained: State of the art language model for NLP, <https://towardsdatascience.com/bert-explained-state-of-the-art-language-model-for-nlp-f8b21a9b6270>
- [6] Alan Akbik, Duncan Blythe, Roland Vollgraf *Contextual String Embeddings for Sequence Labeling*
- [7] Hidayat Ur Rahman, Thomas Hahn, Dr. Richard Segall *Biomedical Disease Name Entity Recognition Using NCBI Corpus*
- [8] <https://towardsdatascience.com/from-pre-trained-word-embeddings-to-pre-trained-language-models-focus-on-bert-343815627598>
- [9] <https://www.sciencedirect.com/science/article/abs/pii/S1386505619310068>
- [10] <https://spacy.io/usage/training>