# Biomedical Named Entity Recognition

ABIDAR Bouchra
BOUSSEBAINE Mustapha
LARBI Abderrahman

Supervised by :
AFFELDT Severine
LABIOD Lazhar

May 29, 2020
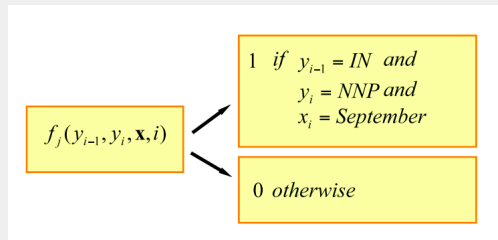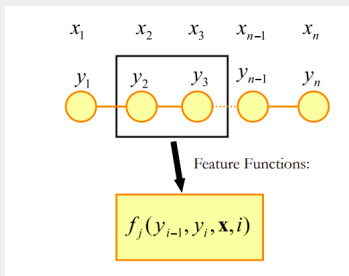
UNIVERSITÉ PARIS DESCARTES
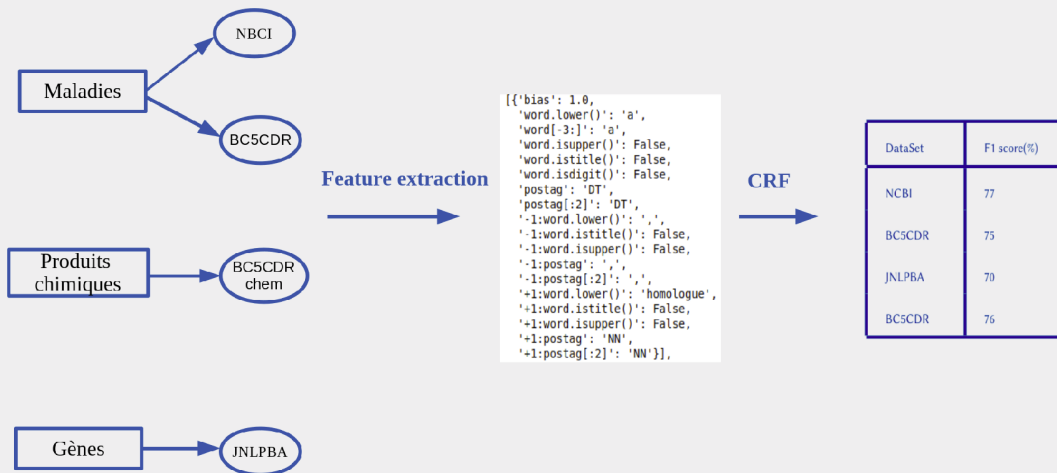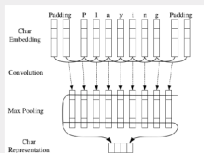
Université de Paris
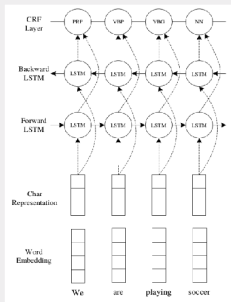
- Discriminant model
- Sequential data prediction
- Features function

La couche CNN



Modéle BiLSTM-CNN-CRF

| DataSet | F1 Score(%) |
|---|---|
| NCBI | 70 |
| BC5CDR | 65 |
| JNLPBA | - |
| BC5CDR-chem | 60 |

**Score de Performance**

- CNN model for generating character embedding.
- Bi-directional LSTM for Word-Level Encoding.
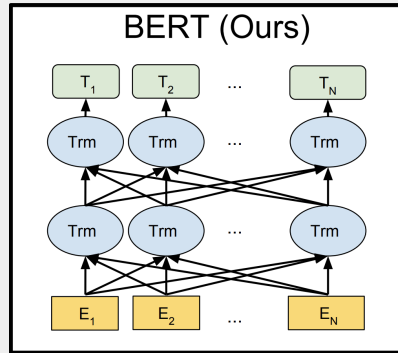- Conditional Random Fields(CRF) for output decoding

- Pre-training deep bidirectional representation models
- Fine tuned with just a single additional layer to produce state of the art results
- Applied on a wide range of NLP STOA Tasks.

- Based on Transformers
- Uses sub word tokenization
  "Here is the sentence I want embeddings for."
  ['[CLS]', 'here', 'is', 'the', 'sentence', 'i', 'want', 'em', '##bed', '##ding', '##s', 'for', '.', '[SEP]']
- Masked LM (MLM)
- Nb : In Named Entity Recognition : The output of the transformer bloc passes by a classification last layer to assign the tag.

Flair Architecture

| Algorithmes | | F1 score(%) |
|---|---|---|
| Flair | NCBI | 88 |
| | BC5CDR | 88 |
| | JNLPBA | 81 |
| | BC5CDR-chem | 93 |

Flair's results

- spaCy is a free, open-source library for advanced Natural Language Processing (NLP) in Python.
- spaCy is designed specifically for production use
- spaCy is not research software.
- spaCy can be used to build information extraction or natural language understanding, or to pre-process text for deep learning.
- spaCy has many features and capabilities.
- spaCy is based on ELMo architecture

ELMo

- in order to train spaCy's NER model we must transform the training data to JSON format. for exemple : ("Who is keyser soze?", "entities": [(7, 17, "PERSON")] )
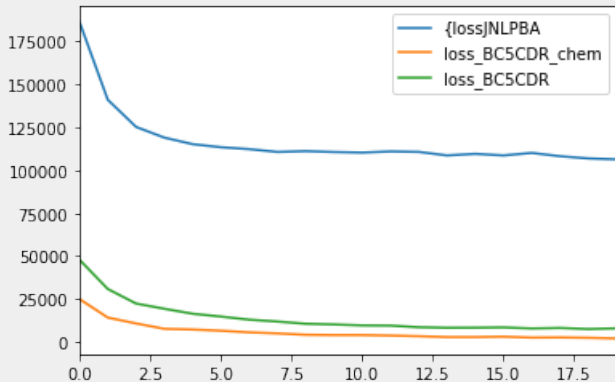


**Figure 1:** spaCy Loss

```
Sentance  :   gene expression and NF-kappa B activation through CD28 requires reactive oxygen pro
duction by 5-lipoxygenase .
......
Tagged  :   gene expression and NF-kappa <B-protein> B <E-protein> activation through CD28 <S-pro
tein> requires reactive oxygen production by 5-lipoxygenase .
--------
Sentance  :   In the in vivo study , the administration ( 50 mg / kg , i . p . ) of TET and FAN i
n mice showed the inhibition of thrombosis by 55 % and 35 % , respectively , while acetylsalicyl
ic 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 .
--------
Sentance  :   OBJECTIVES : The United Kingdom Parkinson ' s Disease Research Group ( UKPDRG ) tri
al found an increased mortality in patients with Parkinson ' s disease ( PD ) randomized to rece
ive 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> Dis
ease <I-Disease> Research Group ( UKPDRG ) trial found an increased mortality in patients with P
arkinson <B-Disease> ' <I-Disease> s <I-Disease> disease <I-Disease> ( PD <B-Disease> ) randomiz
ed to receive 10 mg selegiline <B-Chemical> per day and L <B-Chemical> - <I-Chemical> dopa <I-Ch
emical> compared with those taking L <B-Chemical> - <I-Chemical> dopa <I-Chemical> alone .
```

| 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 |

1- continuous learning (CL)

- The model (CL) takes a tagger as input
- The data entry passes by the three models CRF BERT and SPACY
- Each model yields its results
- The CL model selects the best model on that entry
- The CL starts the retraining with the output dag for the models that yielded the wrong tag.

# Reference

📄 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

📄 Xuezhe Ma and Eduard Hovy *End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF*,2016

📄 Eliyahu Kiperwasser,Yoav Goldberg *Simple and Accurate Dependency Parsing Using Bidirectional LSTM Feature Representations*,2016

📄 DEVLIN, Jacob, CHANG, Ming-Wei, LEE, Kenton, et al. *Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805* , 2018.

THANK YOU!
QUESTIONS?