

NLP 2020

Project presentation

Natural Language Processing A.Y 19/20

Andrea Bacciu

Professor: Roberto Navigli

Teaching assistants:

Edoardo Barba, Niccolò Campolungo, Simone Conia,
Caterina Lacerra, Luigi Procopio

Sapienza University of Rome



Named Entity Recognition - NER

Homework 1



Named Entity Recognition - NER

NER aims to Identify and classify the Named Entities

Classification into 4 classes:

PERson, **ORG**anization, **LOC**ation, **O**ther

John went to California to visit Google

PER **O** **O** **LOC** **O** **O** **ORG**



Task information

Dataset:

- English
- Public
- Metrics: Macro F1-Score
- Splitted in train, dev, test
- Already tokenized

Entity	Train	Dev	Test
O	2,177,423	315,809	335,567
PERSON	100,409	14,396	15,349
LOCATION	84,937	12,359	13,475
ORGANIZATION	61,988	9,043	9,645

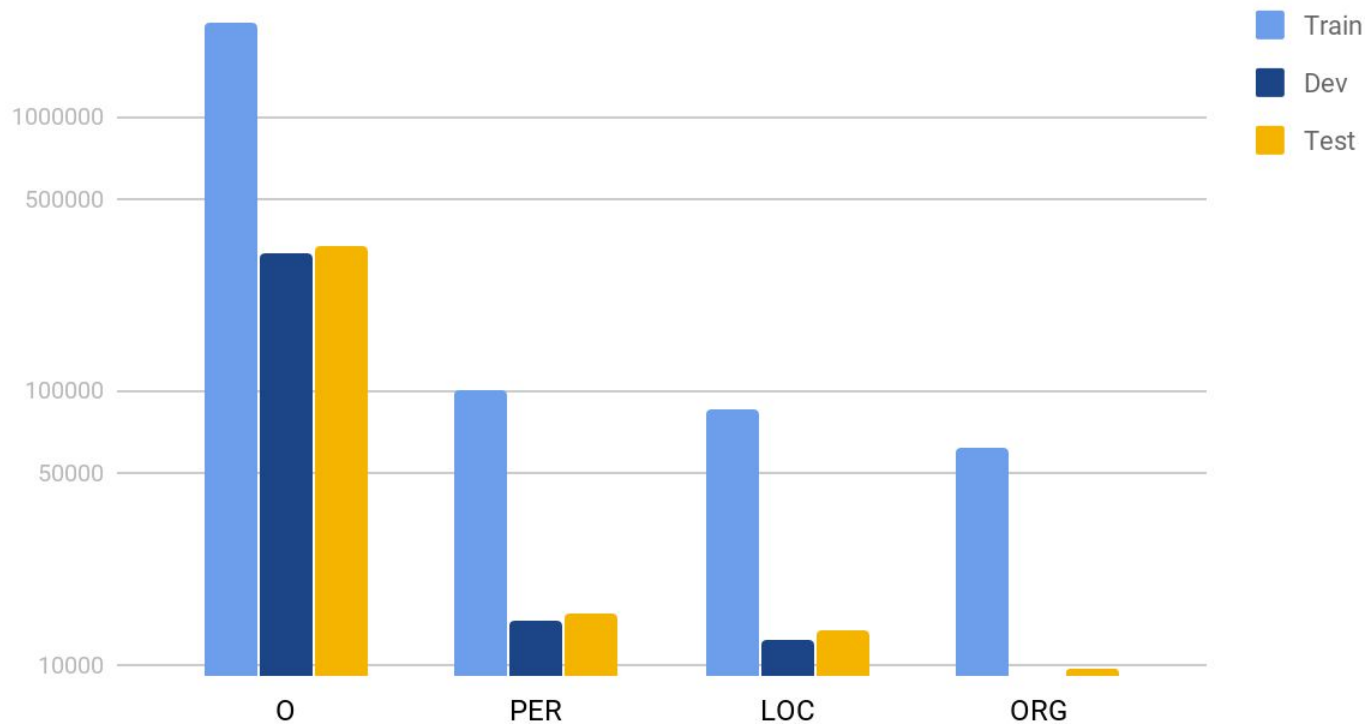
Training:

- Optimizer: Adam
- Learning Rate = 0,001
- Loss function: Cross Entropy
- Epochs: 3
- Batch Size: 32
- DEV F1 early stopping

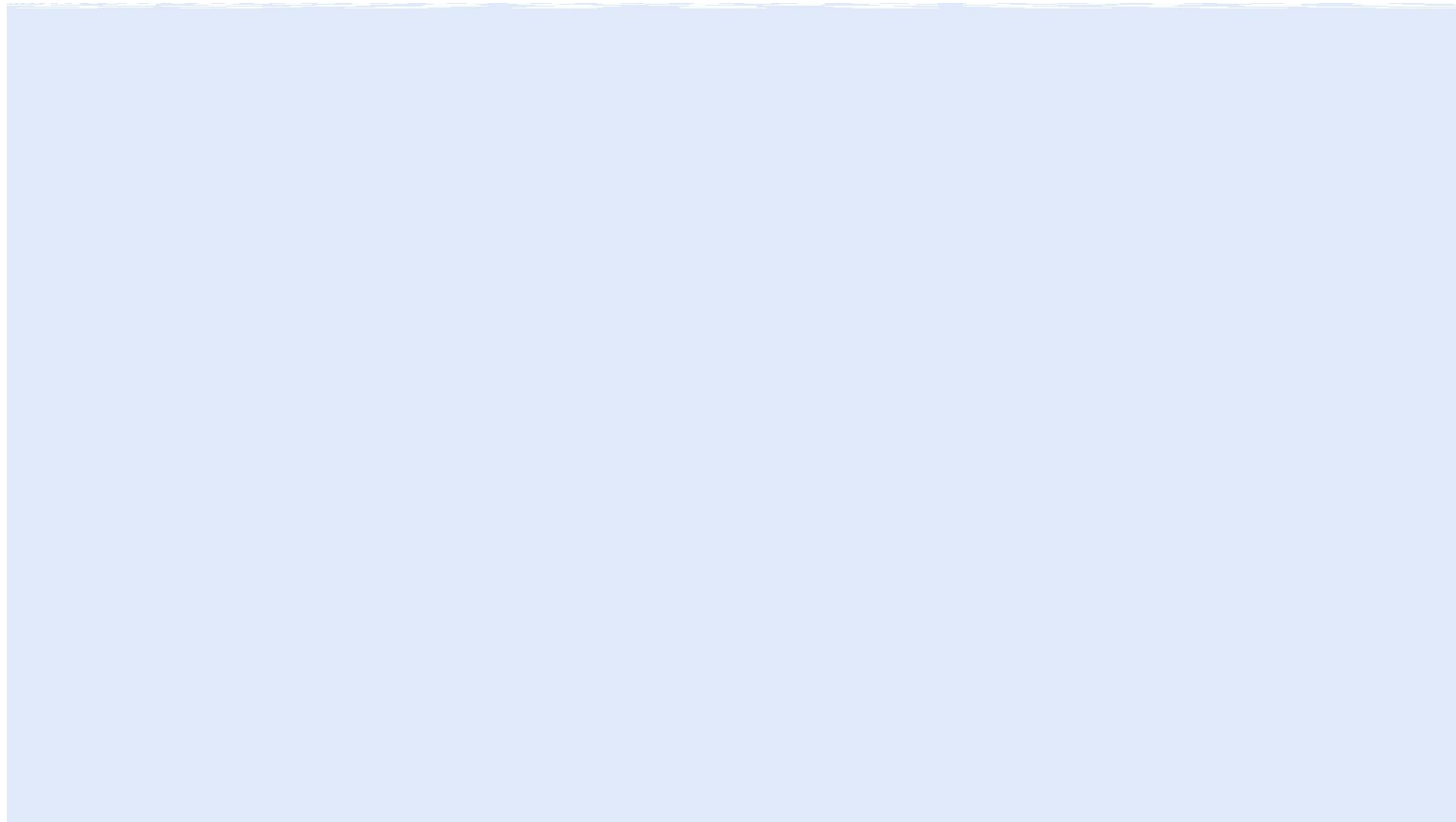


NER Label Distribution

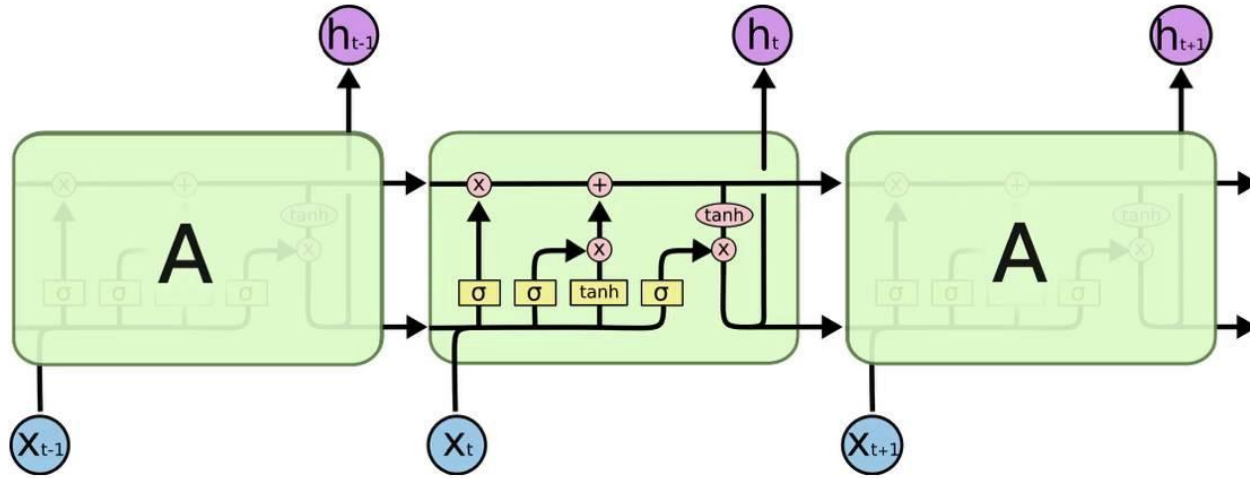
NER Label Distribution - Logarithmic Scaled



Word 2 Vec

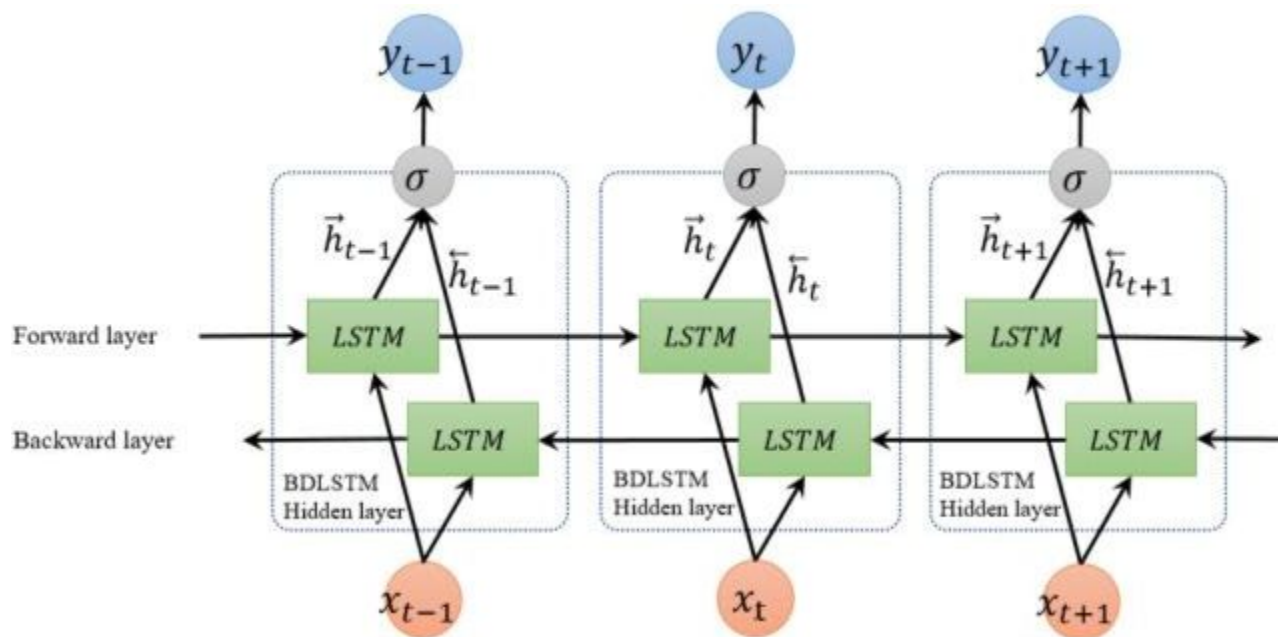


LSTM - Long Short Term Memory



Able to capture long range dependency incorporating memory cell.

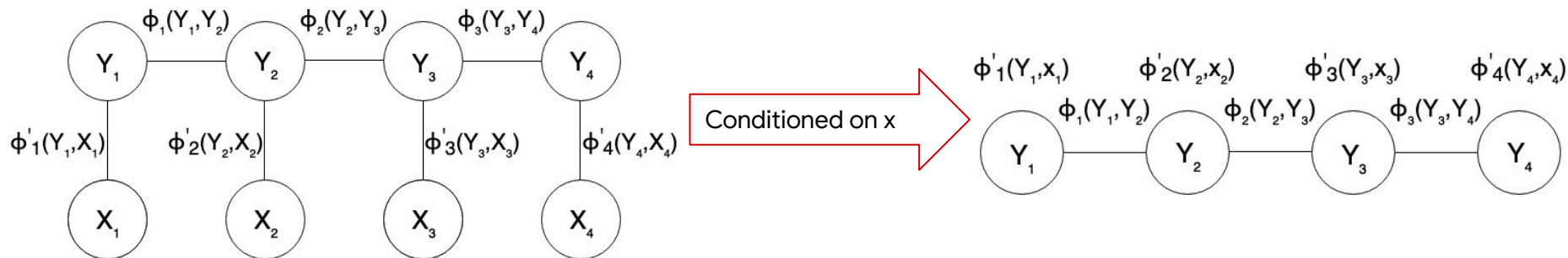
BiLSTM - Bidirectional LSTM



Able to combine the Forward and Backward pass to produce a context-aware output

Conditional Random Field - CRF

CRF structure



CRF helps improve the model's performance **when the current prediction is affected by the previous ones.**

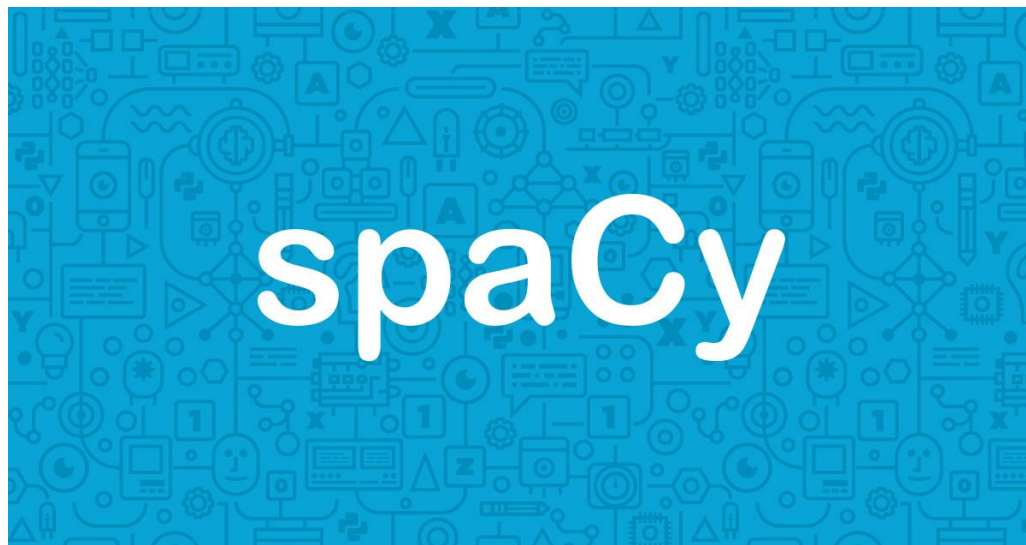
Because the CRF is a **discriminative** model, which are a supervised model class, capable of inferring knowledge from the observed data.

It models the conditional probability $P(Y/X)$ i.e. X is always given or observed. Therefore the graph ultimately reduces to a simple chain.

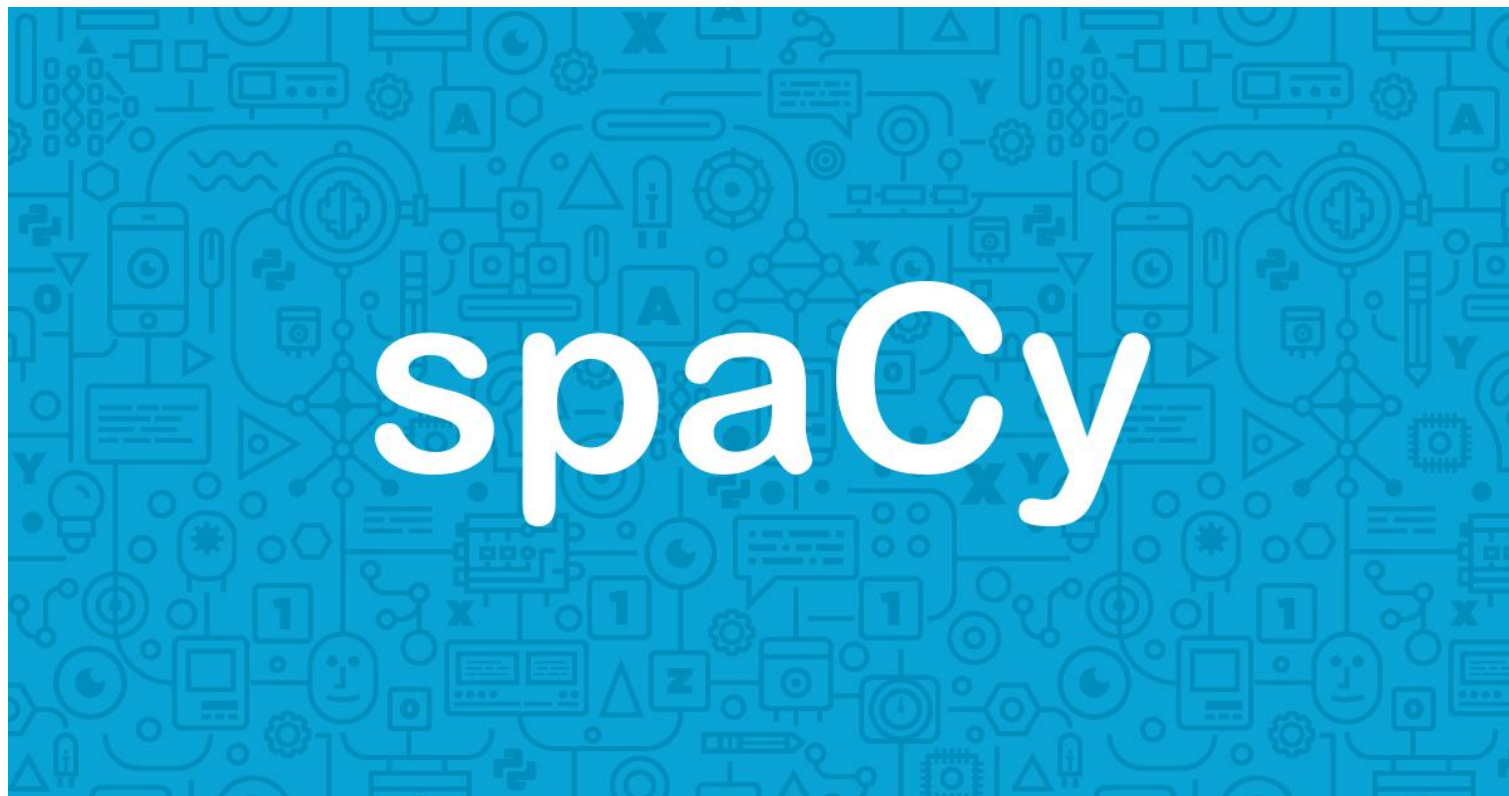


Part of Speech Embeddings

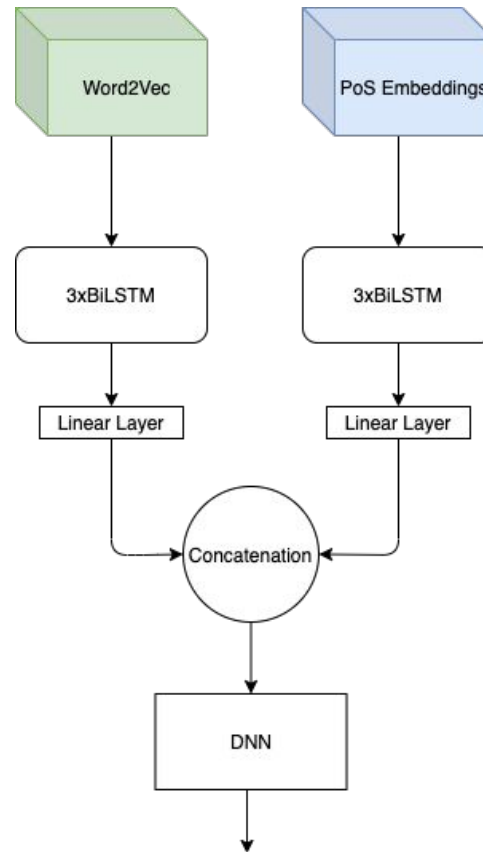
- PoS offer high quality information for the NER task.
- High correlation between labels and PoS tokens.
- I use the SpaCy model to extract the Part of Speech



PoS embeddings with SpaCy



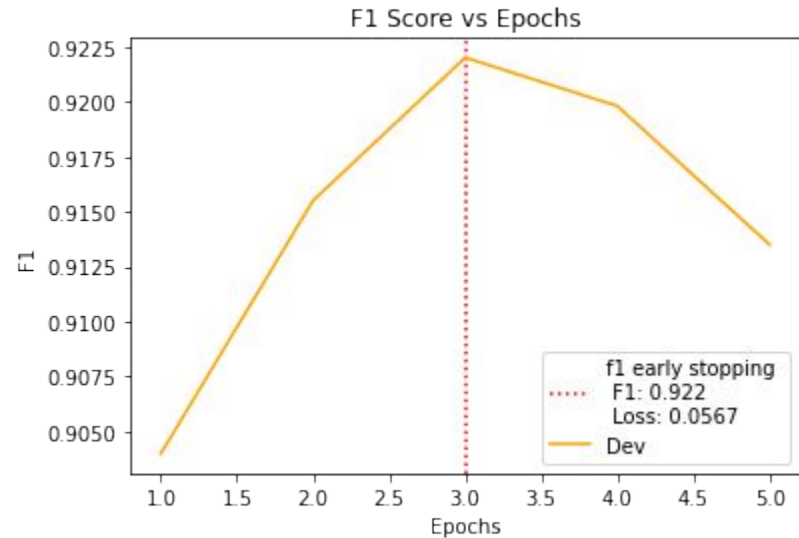
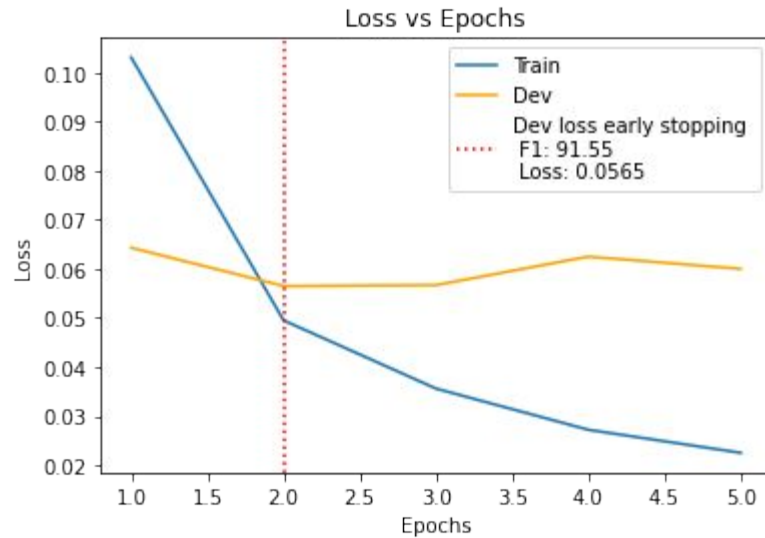
NER Final Model



NER Hyperparameters

HParams	Value	Notes
Epochs	3	Word2Vec
Batch Size	32	
Optimizer	Adam	
Learning Rate	0.001	
Loss Function	Cross Entropy	
Dropout Embeddings	30%	
Word Emb dim	300	
Pos Emb dim	300	
Dropout BiLSTM	30%	
BiLSTM POS	300 out dim	
BiLSTM Words	300 out dim	x3 layer
Linear Layer POS	300	
Linear Layer Words	300	Top Mid Bottom output layer
DNN Linear layer 1	150	
DNN Linear layer 2	75	
DNN Linear layer 3	75	
DNN Linear layer 4	num class	

Early stopping criteria



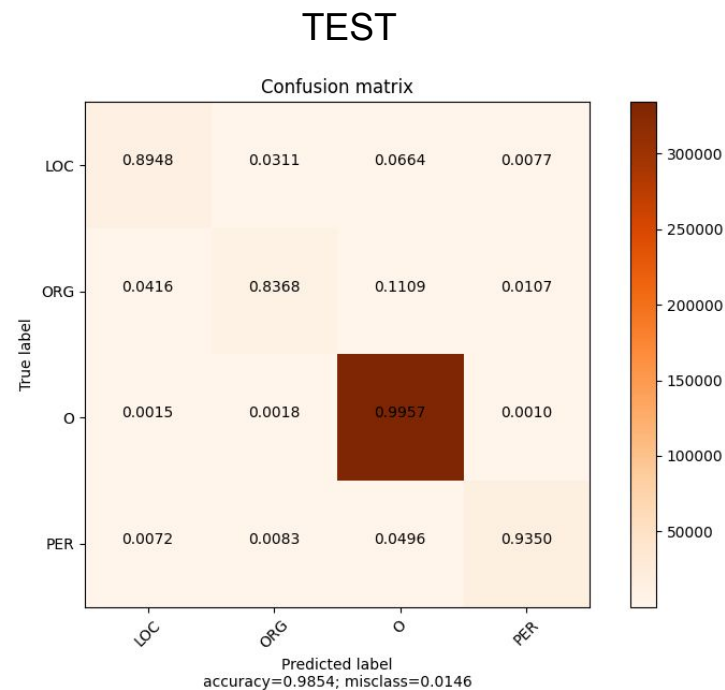
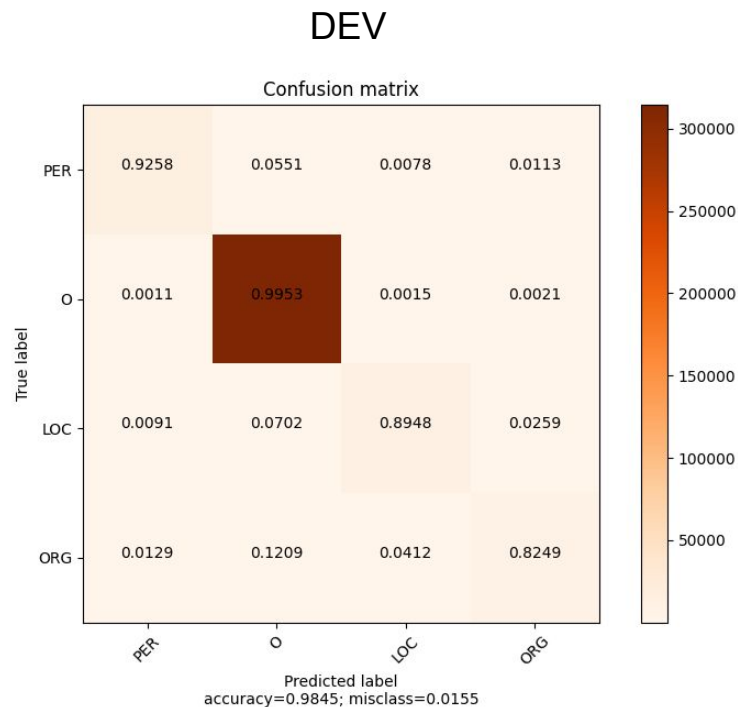
Results

Model	F1 - DEV	F1 - TEST
Baseline	24.88%	24.97%
Word Emb (100d) + BiLSTM	78.02%	78.73%
Word2vec + BiLSTM	80.25%	81.04%
Word2vec + BiLSTM + CRF	88.12%	89.76%
(POS emb + BiLSTM) + (W2V + BiLSTM) + CRF	92.11%	92.37%
(POS emb + BiLSTM) + (W2V + BiLSTM)	92.20%	92.66%

N.B: In the CRF experiments I used the Negative Log Likelihood Loss



Confusion Matrix: Dev vs Test



- The ORG class has worst predictions.
- The O class is better predicted.

Semantic Role Labeling

Homework 2



Semantic Role Labeling (SRL)

SRL is the task of addressing

“Who did What to Whom, How, Where and When?”

the cat ate the fish

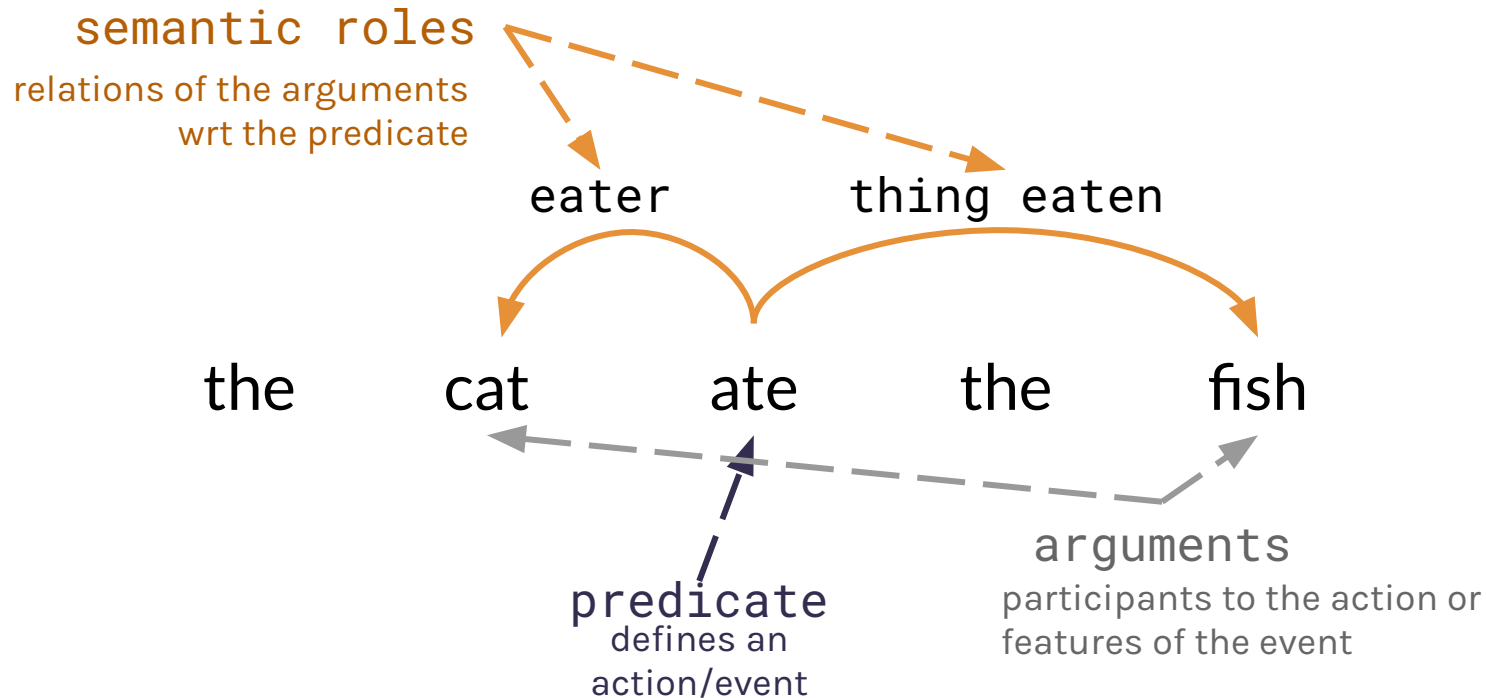
- Who = the cat
- Did what = ate
- Whom = the fish



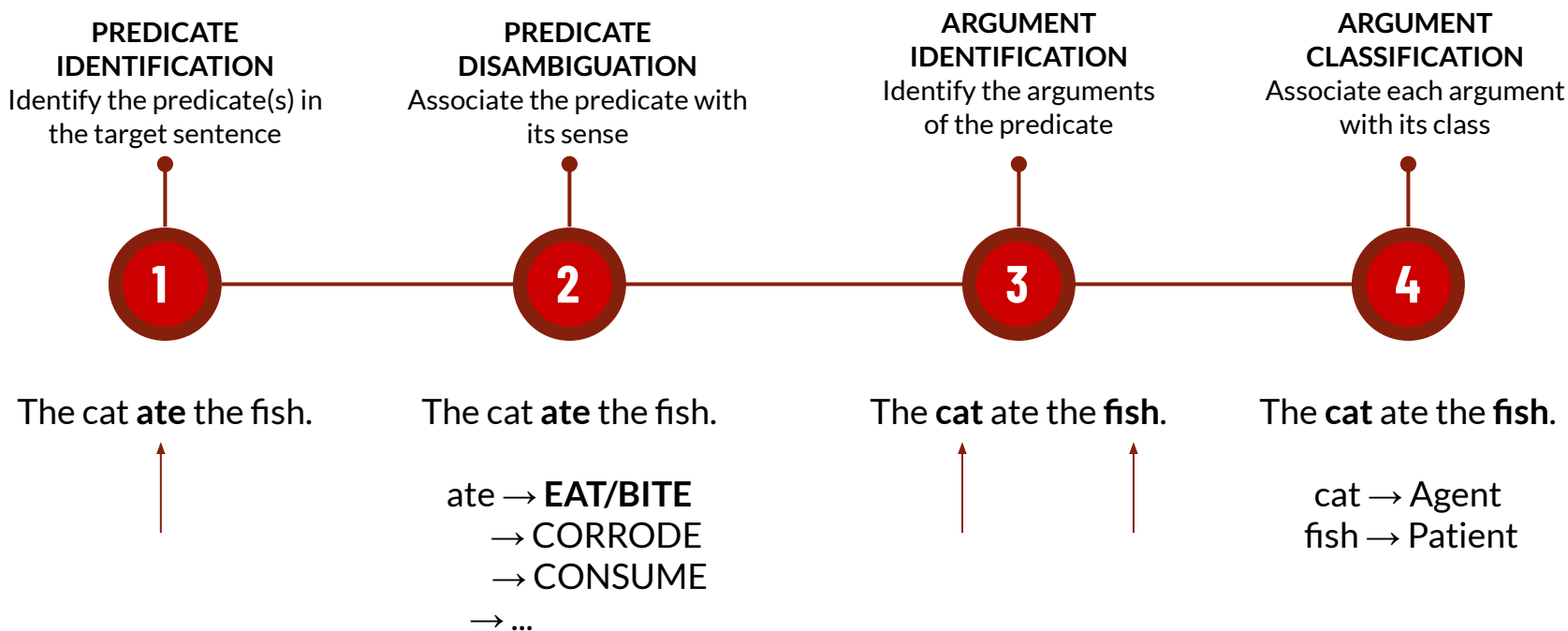
SRL: predicates, arguments and roles

SRL is the task of addressing

“Who did What to Whom, How, Where and When?”



The SRL pipeline



The SRL pipeline

PREDICATE IDENTIFICATION
Identify the predicate(s) in the target sentence



The cat **ate** the fish.



PREDICATE DISAMBIGUATION
Associate the predicate with its sense



The cat **ate** the fish.

ate → **EAT/BITE**
→ CORRODE
→ CONSUME
→ ...

ARGUMENT IDENTIFICATION
Identify the arguments of the predicate



The **cat** ate the **fish**.



ARGUMENT CLASSIFICATION
Associate each argument with its class



The **cat** ate the **fish**.

cat → Agent
fish → Patient



Task information

Dataset:

- English
- Private from Sapienza NLP group
- Metrics: Macro F1-Score
- Split in train, dev, test
- Already tokenized
- Features: PoS, Lemma, Words, Dep. Relations, Dep Heads, Predicates

Training

- Optimizer: Adam
- Learning Rate = 0,001
- Loss function: Cross Entropy
- Epochs: 13
- Batch Size: 32
- Min Frequency = 2 on words and lemmas vocabulary



Sample

sentence_id: {

“words”: [“The”, “cat”, “ate”, “the”, “fish”, “and”, “drank”, “the”, “milk”, “.”],

“lemmas”: [“the”, “cat”, “eat”, “the”, “fish”, “and”, “drink”, “the”, “milk”, “.”],

“pos_tags”: [“DET”, ..., “PUNCT”],

“dependency_relations”: [“NMOD”, ..., “ROOT”, ..., “P”],

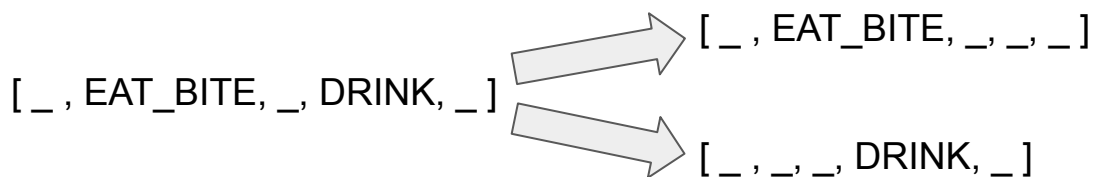
“dependency_heads”: [1, 2, 0, ...],

“predicates”: [“_”, “_”, “EAT_BITE”, “_”, “_”, “_”, “DRINK”, “_”, “_”, “_”],

“roles”: {“2”: [“_”, “Agent”, “_”, “_”, “Patient”, “_”, “_”, “_”, “_”, “_”],

“6”: [“_”, “Agent”, “_”, “_”, “_”, “_”, “_”, “_”, “Patient”, “_”]}

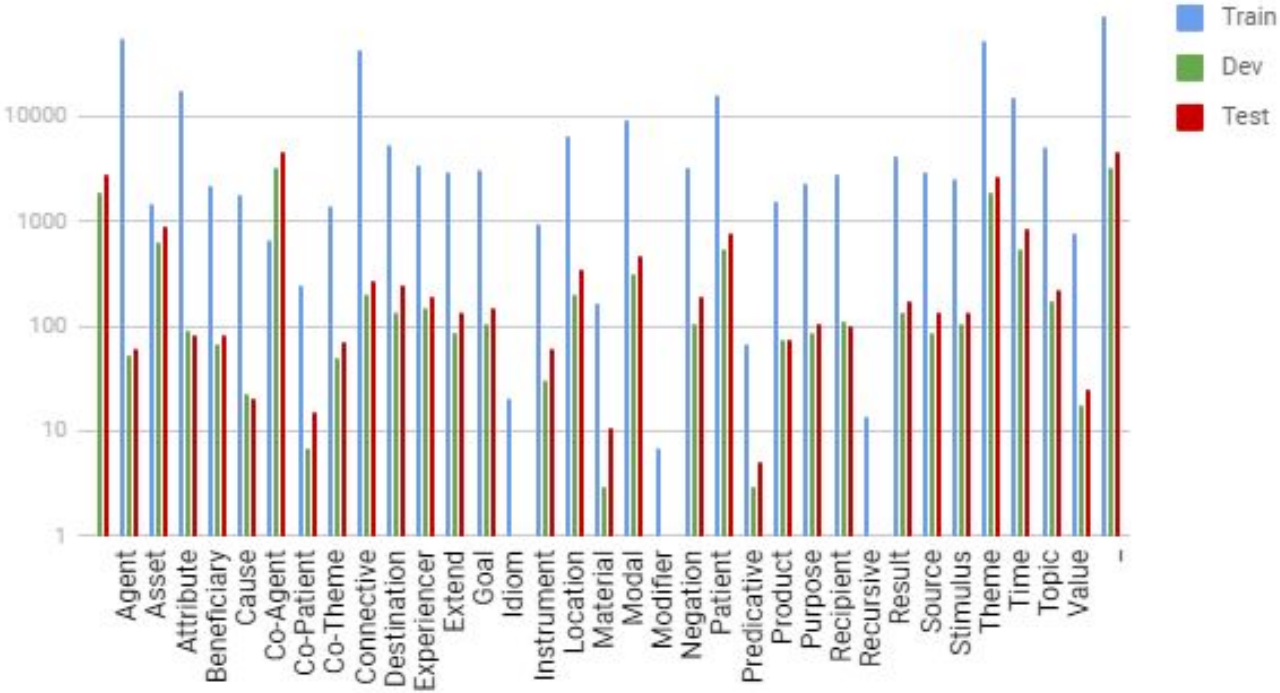
}



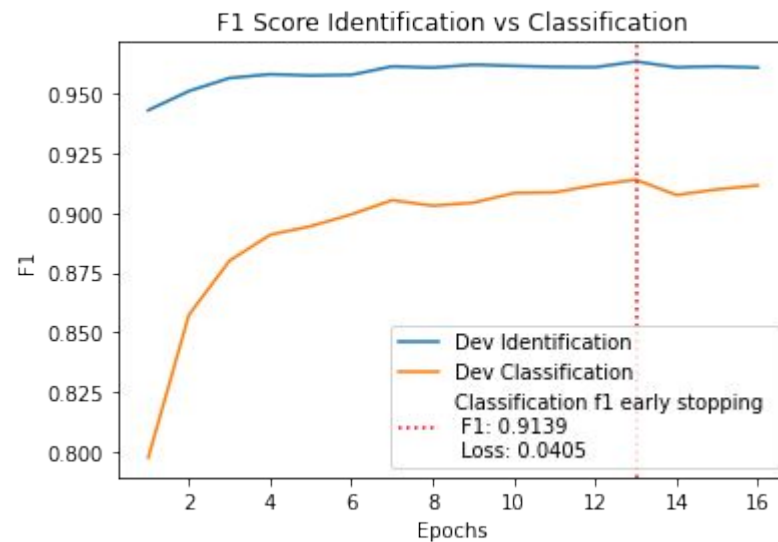
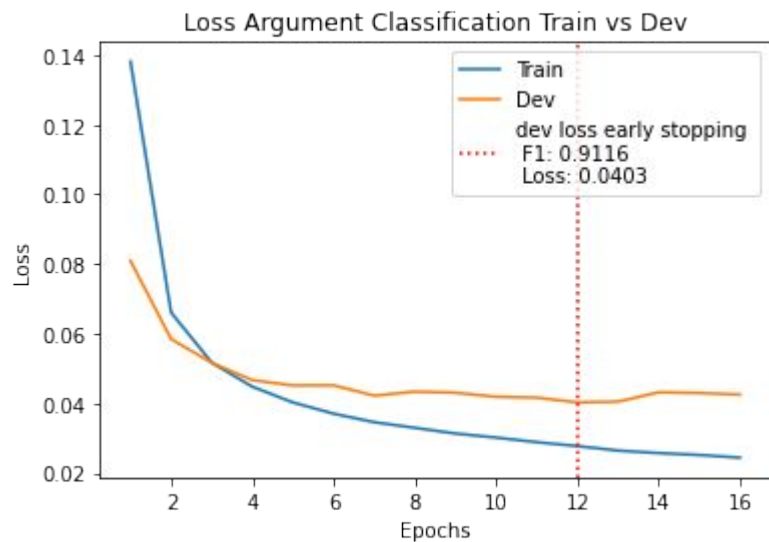
Task information

	DEV	TEST
Idiom	X	X
Modifier	X	X
Recursive	X	✓

Argument label distribution - logarithmic scaled



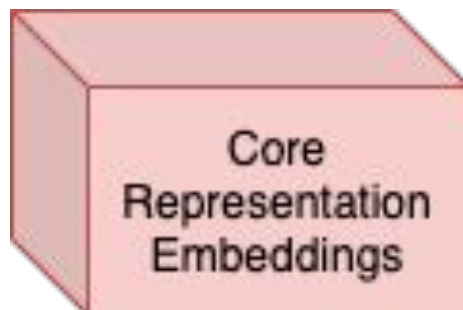
Early stopping criteria



Word Representation

Base WR: Lemma emb \oplus PoS emb \oplus GloVe emb \oplus Predicate emb

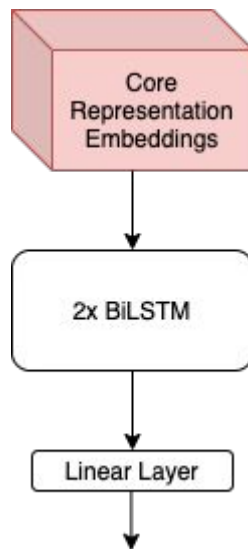
Core WR: Base WR \oplus Dependency Relations Emb



Where \oplus concatenation operator.

The formula are expressed w.r.t a single word in a sentence

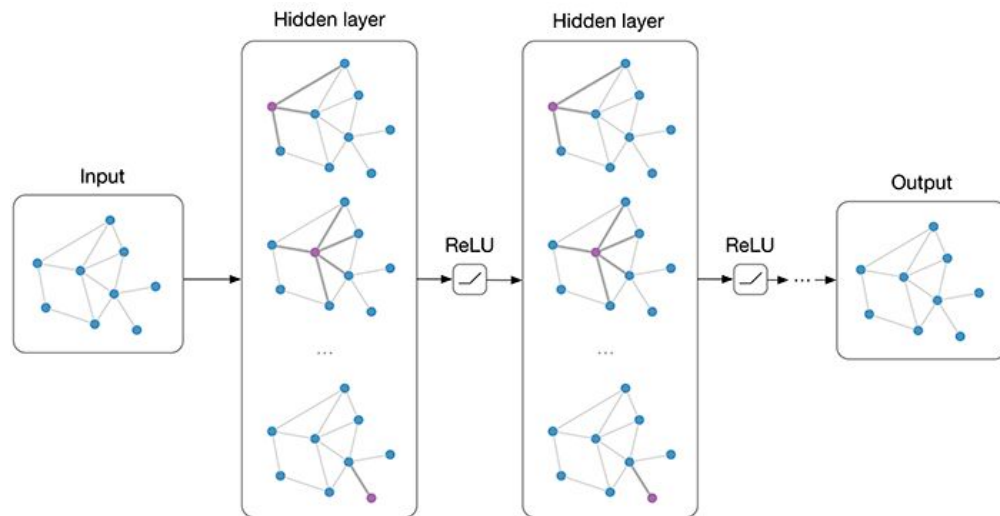
The SRL pipeline



Syntactic Information - Marcheggiani et. al. 2017

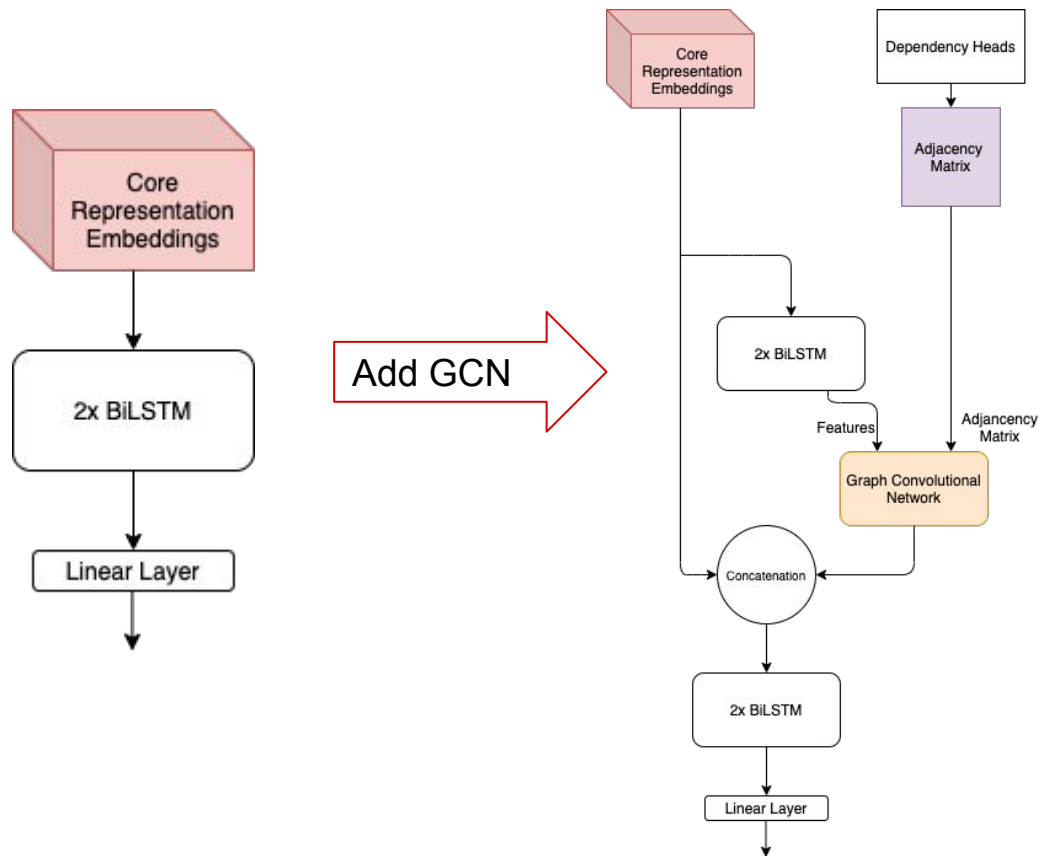
Graph Convolutional Network - (Kipf and Welling)

- GCNs are neural networks that operate on graphs and induce features of nodes based on the properties of their neighbours
- Use Normalized Adjacency Matrix
 - $A * D^{-1}$
- GCN is not able to capture long dependencies between distant nodes in the graph.
 - Solve using context-aware input

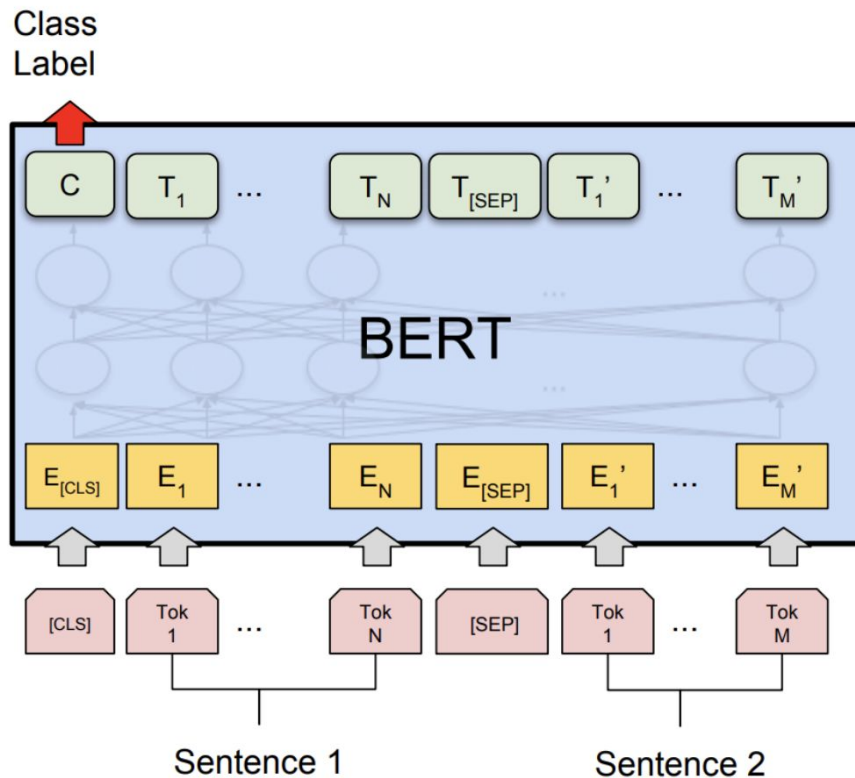


“The semantic representations are closely related to syntactic ones”

The SRL pipeline



BERT - Bidirectional Encoder Representation from Transformer



BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

(Bert Base uncased 768-dim)



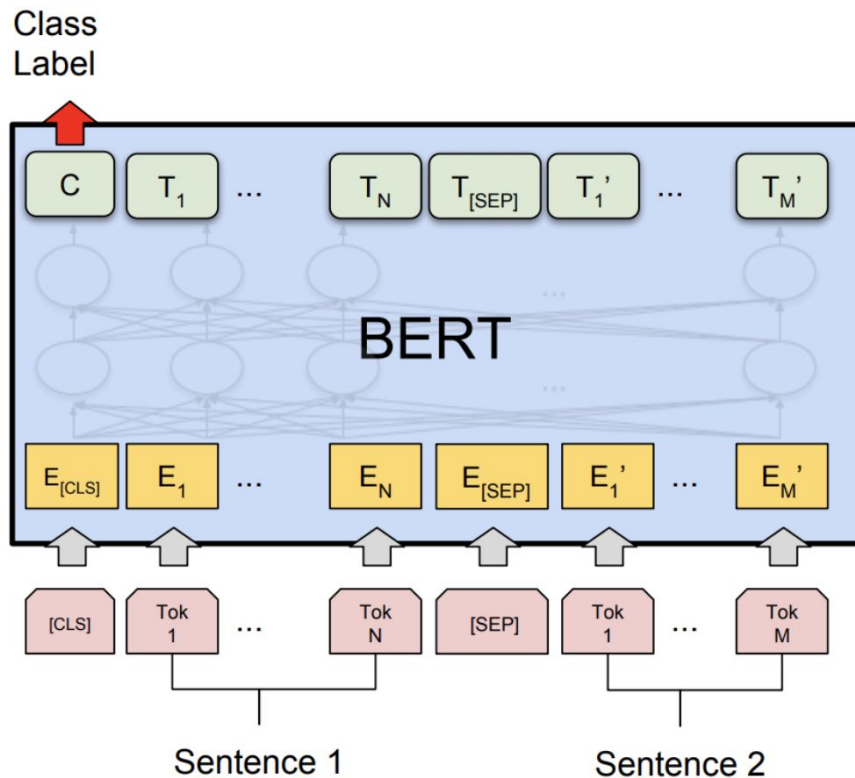
The Non-trainable Bert Embeddings

Vector representation based on the context

Use Word pieces vocabulary

- 30.000 word pieces in vocabulary
- Word segmentation based purely on frequency
- Least frequent words are divided into several word pieces
- *The infrequent 'Embeddings' -> ['em', '##bed', '##ding', '##s']*
- Almost impossible encounter OOV word

I don't apply the fine-tuning to the bert model.



(Bert Base uncased 768-dim)

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

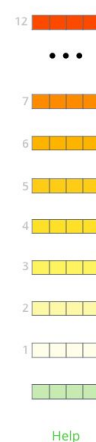
Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

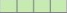

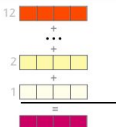

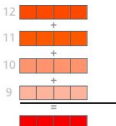



The Non-trainable Bert Embeddings

- To produce the embeddings: I summed the last 4 hidden layers
 - this pooling strategy is proven to be one of the most efficient with low memory consumption
- I removed the special tokens like [CLS] and [SEP]
- I merged the word pieces produced by a single words using the average of their embeddings.

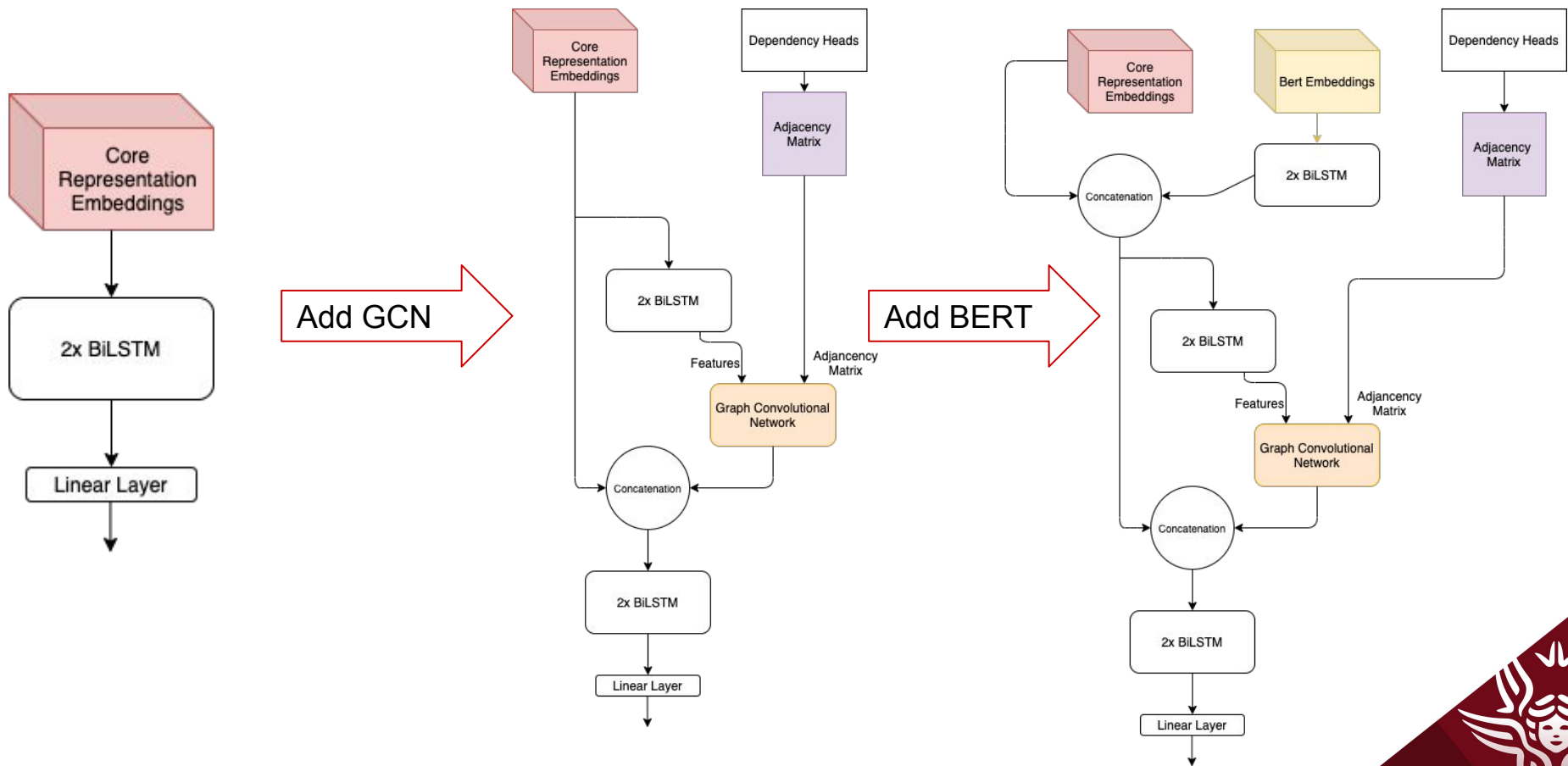
What is the best contextualized embedding for "Help" in that context?
For named-entity recognition task CoNLL-2003 NER



		Dev F1 Score
First Layer	Embedding 	91.0
Last Hidden Layer	12 	94.9
Sum All 12 Layers		95.5
Second-to-Last Hidden Layer	11 	95.6
Sum Last Four Hidden		95.9
Concat Last Four Hidden		96.1

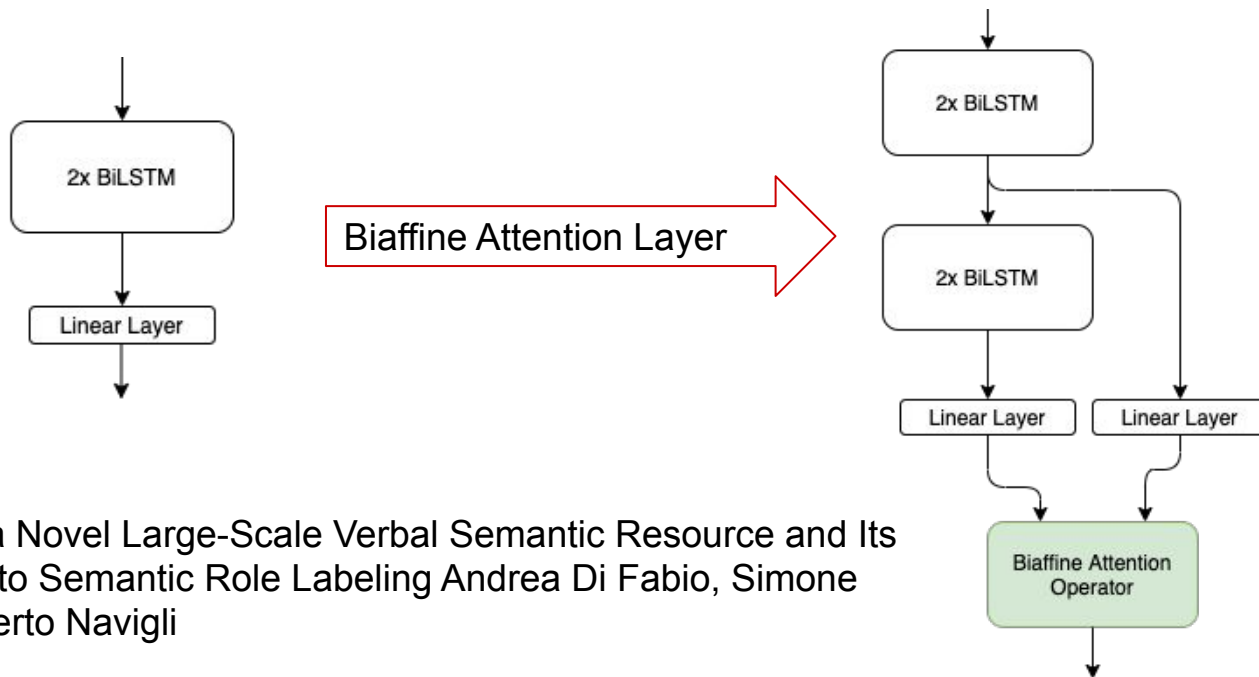


The SRL pipeline



Biaffine Attention Layer

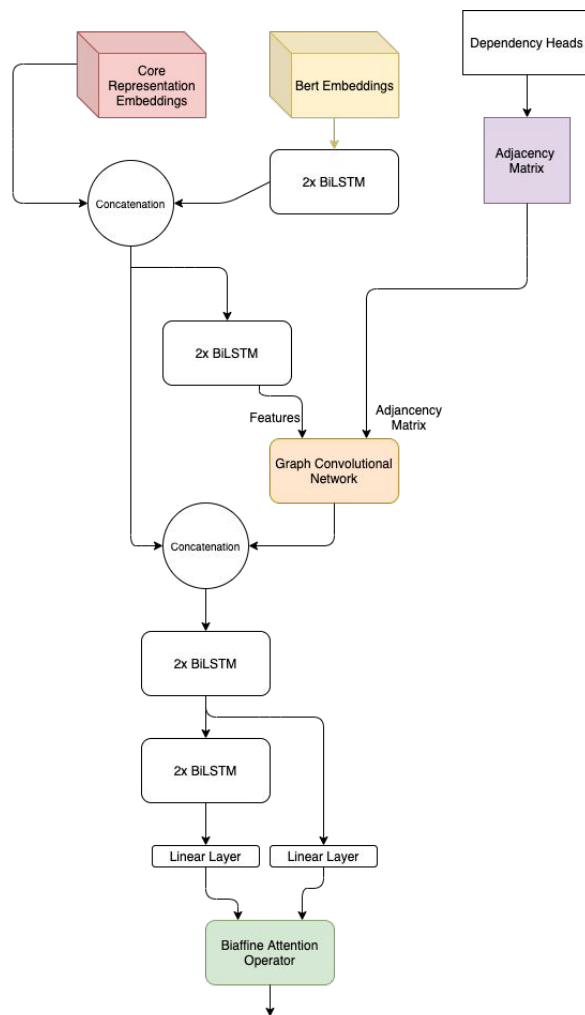
$$\text{Biaffine}(y_1, y_2) = \underbrace{y_1^T \mathbf{U} y_2}_{\text{Bilinear}} + \underbrace{\mathbf{W}(y_1 \circ y_2) + \mathbf{b}}_{\text{Linear}}$$



VerbAtlas: a Novel Large-Scale Verbal Semantic Resource and Its Application to Semantic Role Labeling
Andrea Di Fabio, Simone Conia, Roberto Navigli



Final Model



SRL Hyperparameters

HParams	Value	Notes
Epochs	13	
Batch Size	32	
Optimizer	Adam	
Learning Rate	0.001	
Loss Function	Cross Entropy	
Word Vocab Min frequency	2	
Lemma Vocab Min frequency	2	
Dropout Embeddings	30%	
Word Emb dim	300	GloVe 6B Bert Base uncased
Bert Emb dim	768	
Pos Emb dim	300	
Lemma Emb dim	300	
Predicate Emb dim	400	
Dependency Relations Emb dim	300	
Num GCN-Convolutional Layer	2	
GCN-Convolutional Layer Hidden size	250	first layer
GCN-Convolutional Layer Hidden size	35	second layer
GCN Dropout	50%	
Dropout BiLSTM	30%	
BiLSTM Bert	300 out dim	x2 layer
BiLSTM GCN	300 out dim	x2 layer
BiLSTM Biaffine B1	300 out dim	x2 layer
BiLSTM Biaffine B2	300 out dim	x2 layer
Biaffine inputs	35, 35	y1, y2
Biaffine outputs	35	length of label vocab
Linear layer y1	35	
Linear layer y2	35	



Results

SRL Results				
Experiments	F1-Dev		F1-Test	
	Identification	Classification	Identification	Classification
Baseline + Base WR	90.84%	85.84%	91.80%	87.38%
Baseline + Core WR	94.05%	88.96%	94.52%	90.30%
Baseline + GCN + Core WR	95.68%	89.35%	96.62%	91.25%
Baseline + Bert + GCN + Core WR	95.80%	91.07%	95.90%	92.05%
Baseline + Bert + GCN + Biaffine + Core WR*	96.32%	91.38%	96.87%	92.78%
Baseline + Bert + Core WR	95.14%	90.71%	95.72%	91.98%
Baseline + GCN + Biaffine + Core WR	95.78%	89.70%	96.72%	92.28%
Baseline + Bert + Biaffine + Core WR	95.09%	90.80%	95.90%	92.21%



DEV



TEST



- The _ class is better predicted on both dataset.
- The material class is poorly predicted in both dataset.
- Recursive on the test is poorly predicted, the train dataset has 14 Recursive examples.

Reference

- VerbAtlas: a Novel Large-Scale Verbal Semantic Resource and Its Application to Semantic Role Labeling Andrea Di Fabio, Simone Conia, Roberto Navigli
VerbAtlas: a Novel Large-Scale Verbal Semantic Resource and Its Application to Semantic Role Labeling Andrea Di Fabio, Simone Conia, Roberto Navigli
- Encoding Sentences with Graph Convolutional Networks for Semantic Role Labeling Diego Marcheggiani, Ivan Titov
- Semi-supervised classification with graph convolutional networks Thomas N. Kipf, Max Welling
- BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
- Pytorch-crf package



Thank you for the attention

Andrea Bacciu

