# WEB CONTENT MINING
# NAMED ENTITY RECOGNITION

## Elina Tugaeva, Eliana Ruslanova, Maida Nazifi, Nikolaj Kolbasko, Samail Guliev

Web Mining IE684 Course
Faculty of Business Informatics and Mathematics
University of Mannheim, 2020

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. LSTM-CRF
5. RNN - LSTM
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

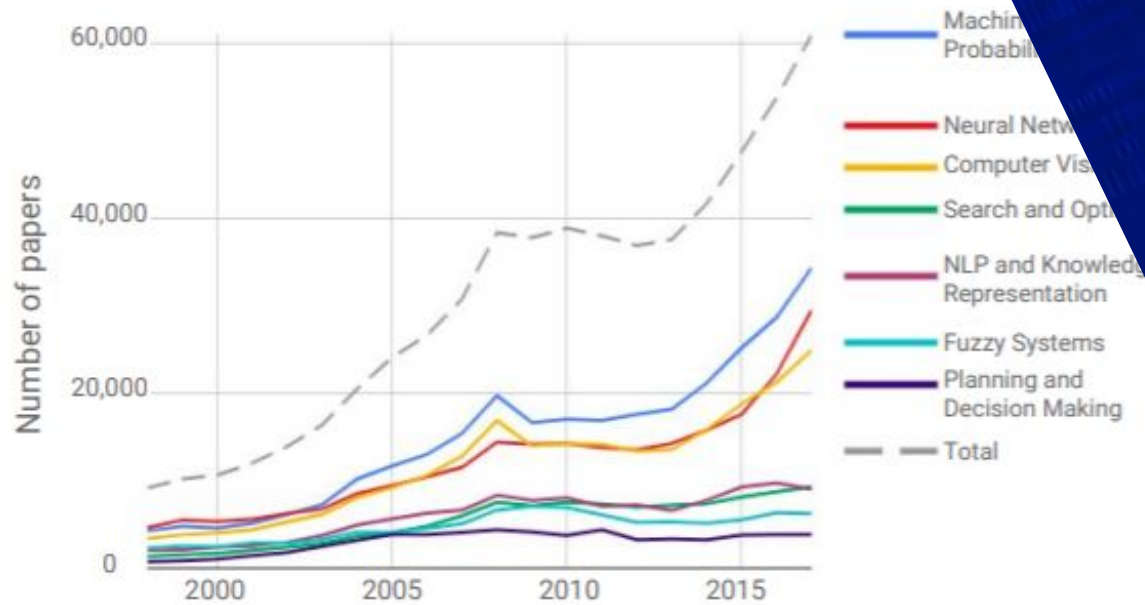# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. LSTM-CRF
5. RNN - LSTM
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# MOTIVATION

Number of AI papers on Scopus by subcategory (1998—2017)
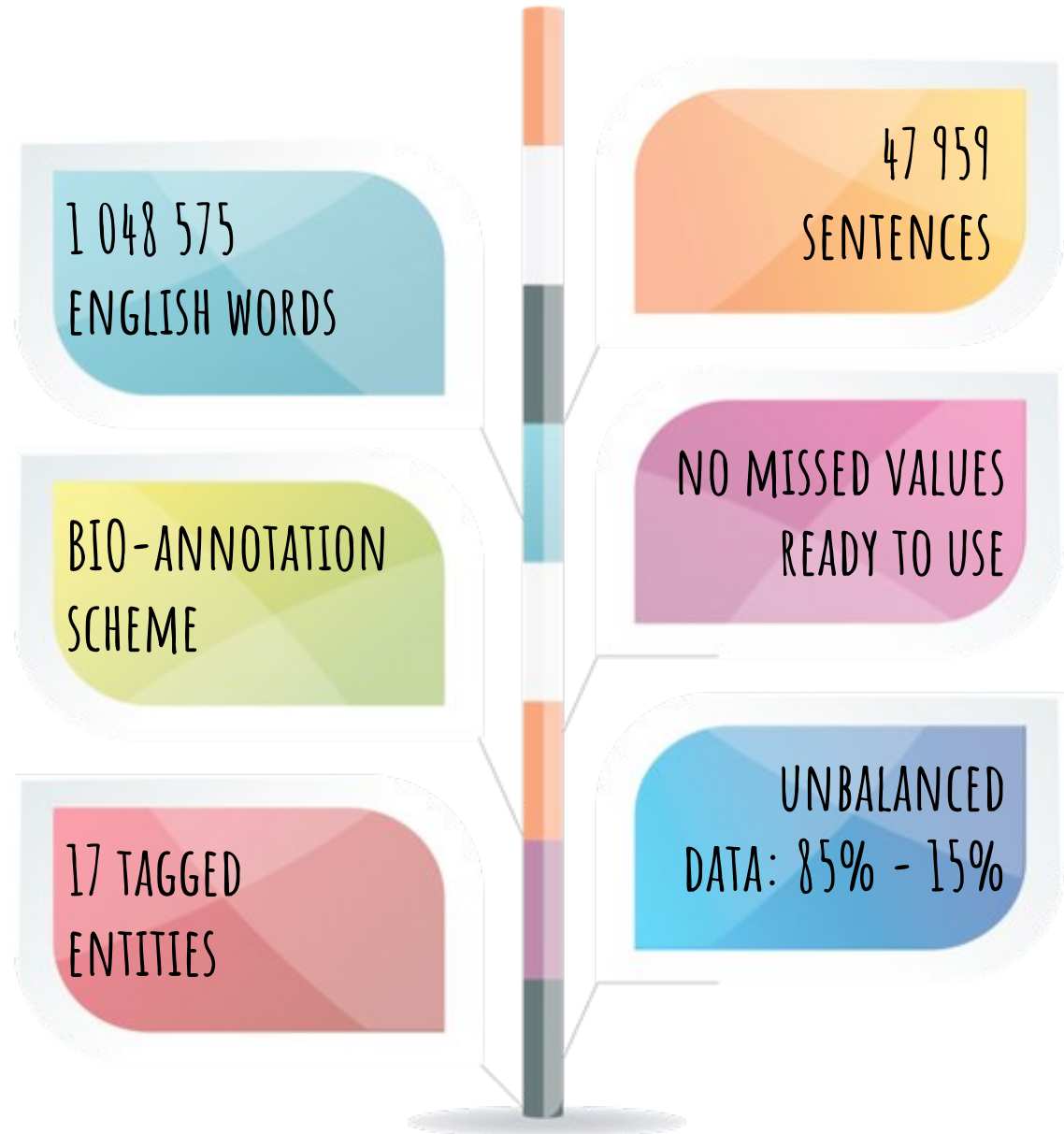Source: Elsevier

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. LSTM-CRF
5. RNN - LSTM
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# Data Set

Annotated Corpus for

Named Entity Recognition

from Kaggle

| | Sentence # | Word | POS | Tag |
|---|---|---|---|---|
| 0 | Sentence: 1 | Thousands | NNS | O |
| 1 | NaN | of | IN | O |
| 2 | NaN | demonstrators | NNS | O |
| 3 | NaN | have | VBP | O |
| 4 | NaN | marched | VBN | O |
| 5 | NaN | through | IN | O |
| 6 | NaN | London | NNP | B-geo |

1 048 575 ENGLISH WORDS

BIO-ANNOTATION SCHEME

17 TAGGED ENTITIES

47 959 SENTENCES

NO MISSED VALUES READY TO USE

UNBALANCED DATA: 85% - 15%

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. LSTM-CRF
5. RNN - LSTM
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# Conditional random field

- A probabilistic, discriminative classification model for sequences
- Directly model the association between the observed features and labels for those features: defines a posterior probability of a label sequence given an input observation sequence

$$p(s_1, .., s_m | x_1, .., x_m, w) = \frac{exp(w\Phi(x,s))}{\sum exp(w\Phi(x,s))}$$

$$L = \sum_{i=1}^{n} p(s^i | x^i, w) - \frac{\lambda_2}{2} \times |w|^2 - \lambda_1 \times |w|$$
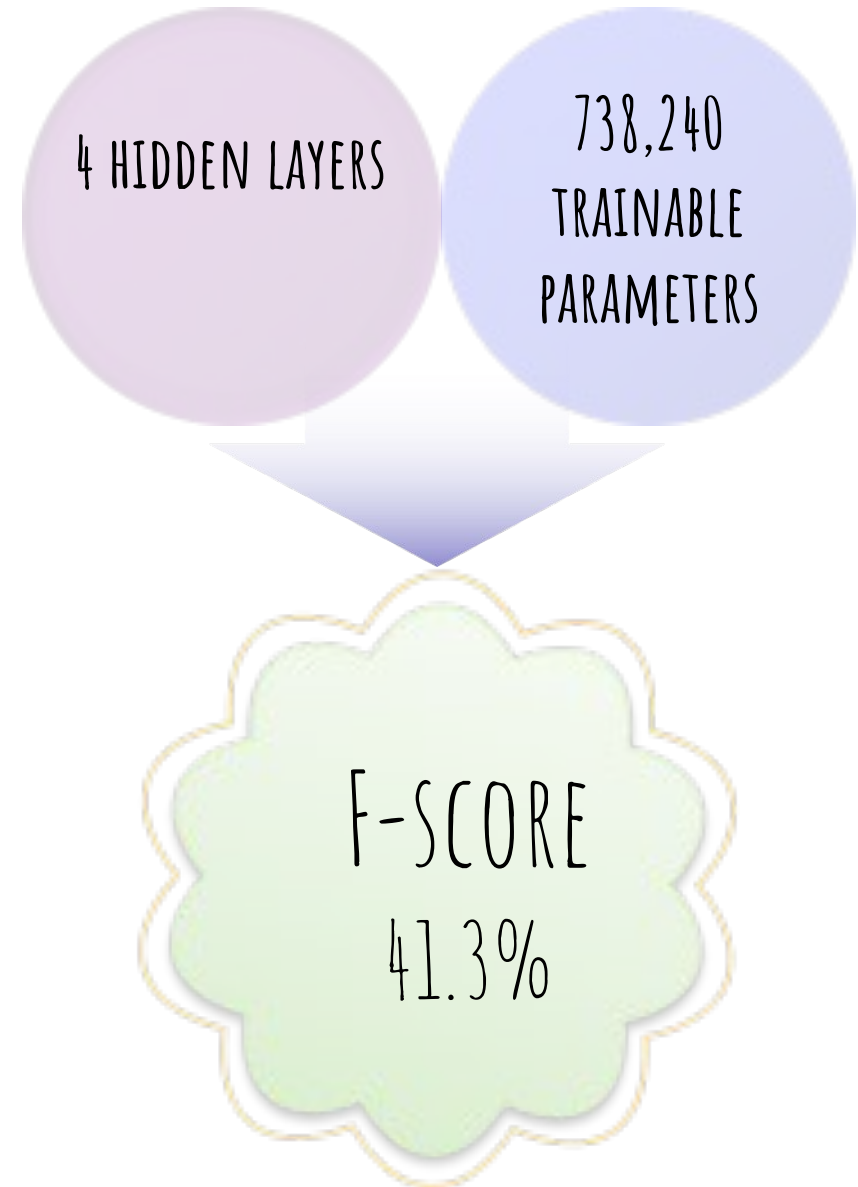
REGULARIZED

F-SCORE

90%

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. RNN - LSTM
5. LSTM-CRF
6. Character embeddings
7. LSTM-CNN-CRF
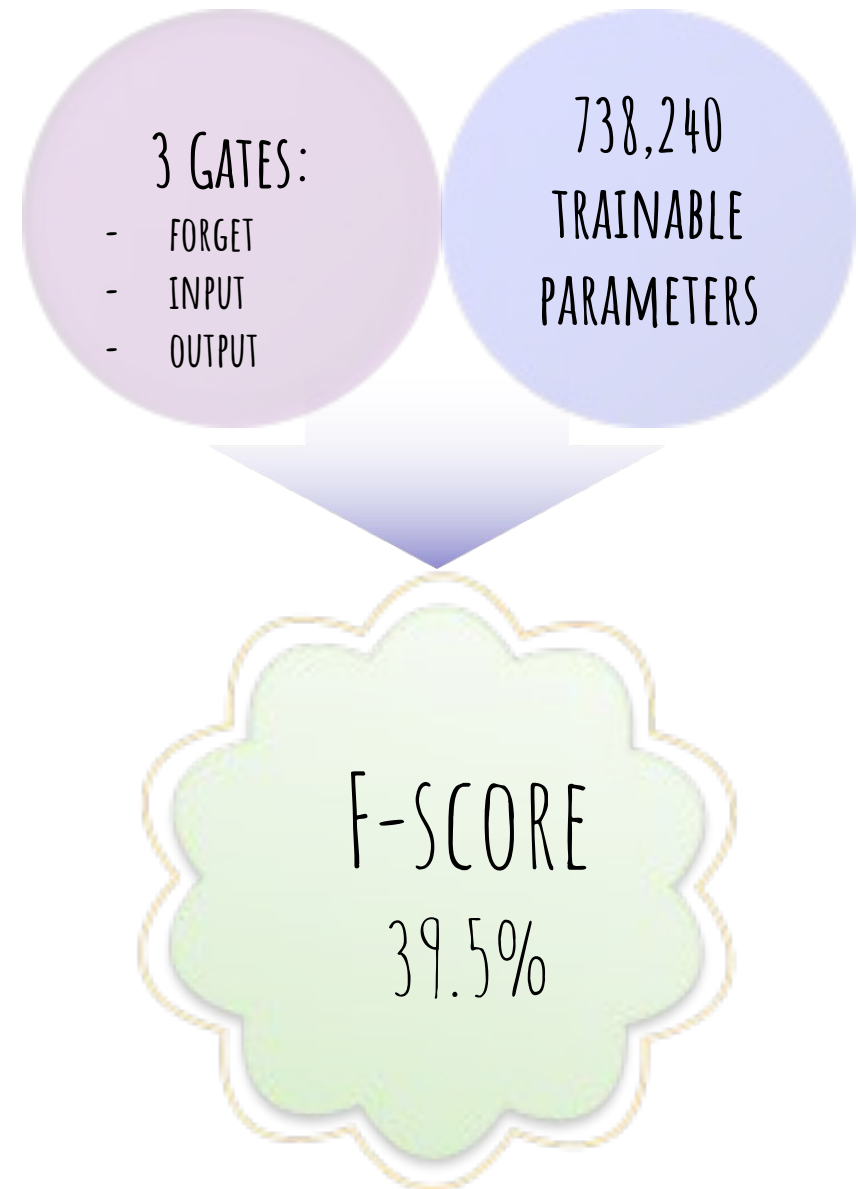8. BERT
9. SpaCy

# RNN

- Deep learning Method known for good results in NER Tasks
  - BiDirectional Rnn

- Embedding layer fed into Rnn with a Dropout Layer = 0.5

- long-term dependency problem

4 hidden layers

738,240 trainable parameters

F-score
41.3%

# Lstm

- Trains model using back-propagation

- tanh activation function

- Even though not good general performance, results in geo and per tags
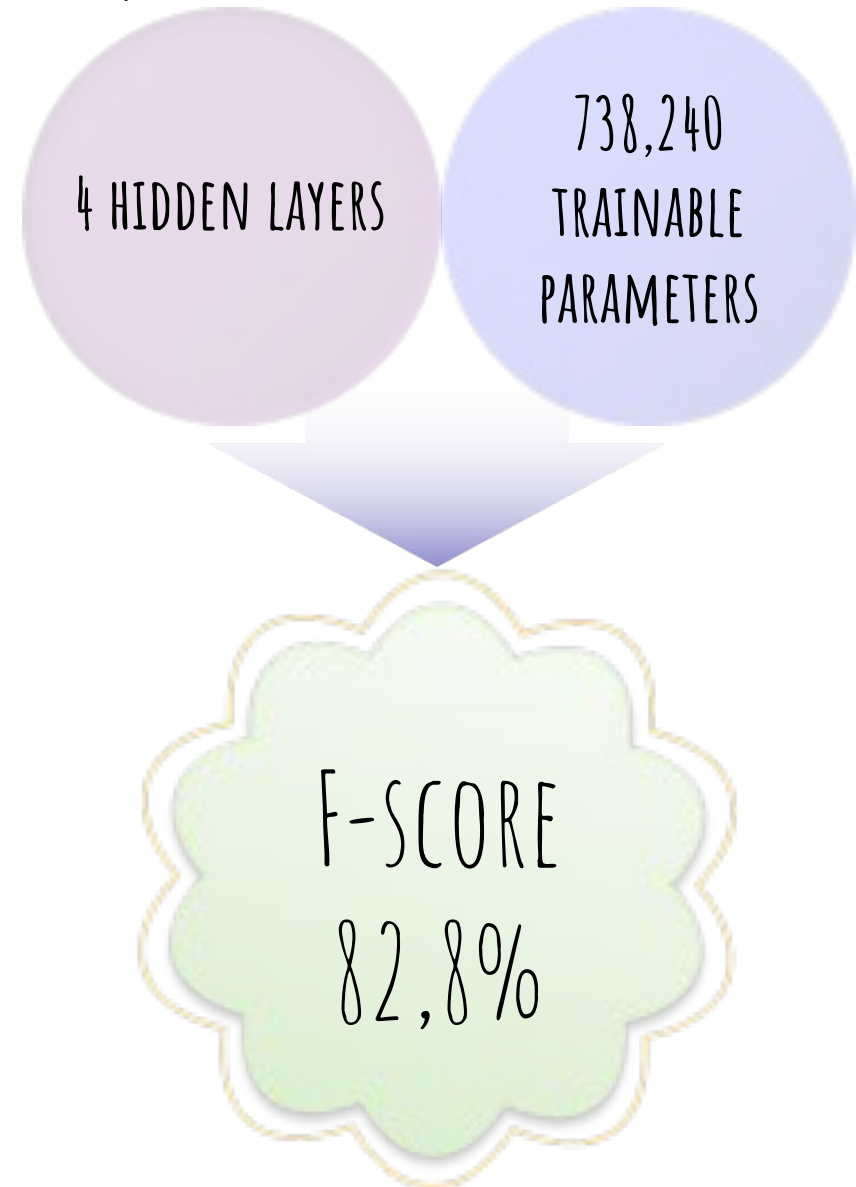
- combination with crf or cnn

3 Gates:
- forget
- input
- output

738,240 trainable parameters

F-score
39.5%

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. RNN - LSTM
5. LSTM-CRF
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# Long term short memory CRF

- CRF built on top of a BI-LSTM:
    CRF layer is basically an optimisation on top of BI-LSTM layer
- Embedding layer in a vector form is passed to the crf algorithm:
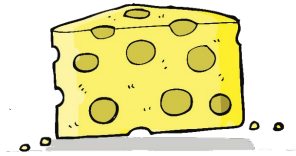    => little effort on feature engineering

4 hidden layers

738,240 trainable parameters

F-score 82,8%

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. RNN - LSTM
5. LSTM-CRF
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# CHARACTER EMBEDDINGS
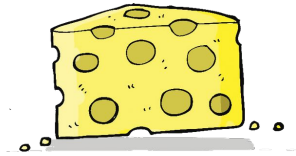
CHEESE

| C | H | E | E | S | E |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 3 | 4 | 3 |

| 1 | 2 | 3 | 3 | 4 | 3 |
|---|---|---|---|---|---|

| 1 | 2 | 3 | 3 | 4 | 3 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

Max size of word vector =10

# CHARACTER EMBEDDINGS

CHEESE

| C | H | E | E | S | E |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 3 | 4 | 3 |

| 1 | 2 | 3 | 3 | 4 | 3 |
|---|---|---|---|---|---|

| 1 | 2 | 3 | 3 | 4 |
|---|---|---|---|---|

,

| 3 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|

MAX SIZE OF WORD VECTOR = 5

# LSTM for char embeddings

WORD EMBEDDINGS

LSTM for char EMBEDDINGS

BI-LSTM

F-score 82,3%

Time for one epoch: 120 seconds

# CNN for char embeddings

WORD EMBEDDINGS

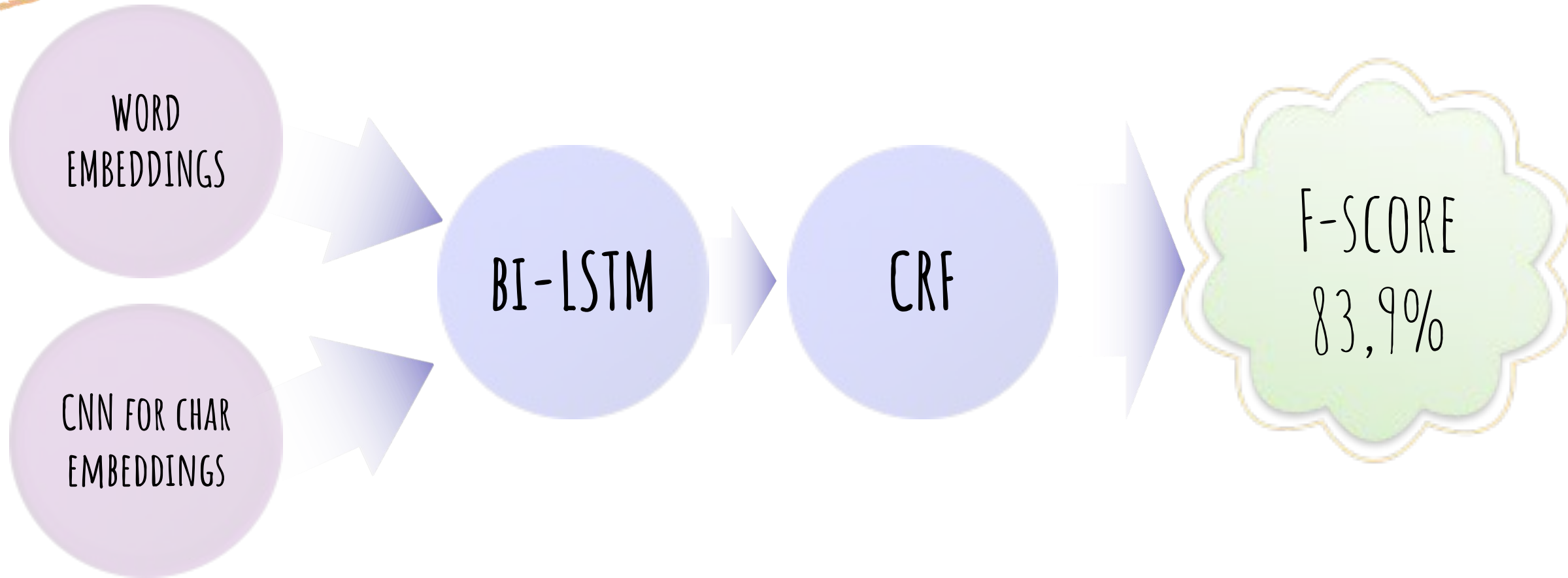CNN for char EMBEDDINGS

BI-LSTM

F-score 79,8%
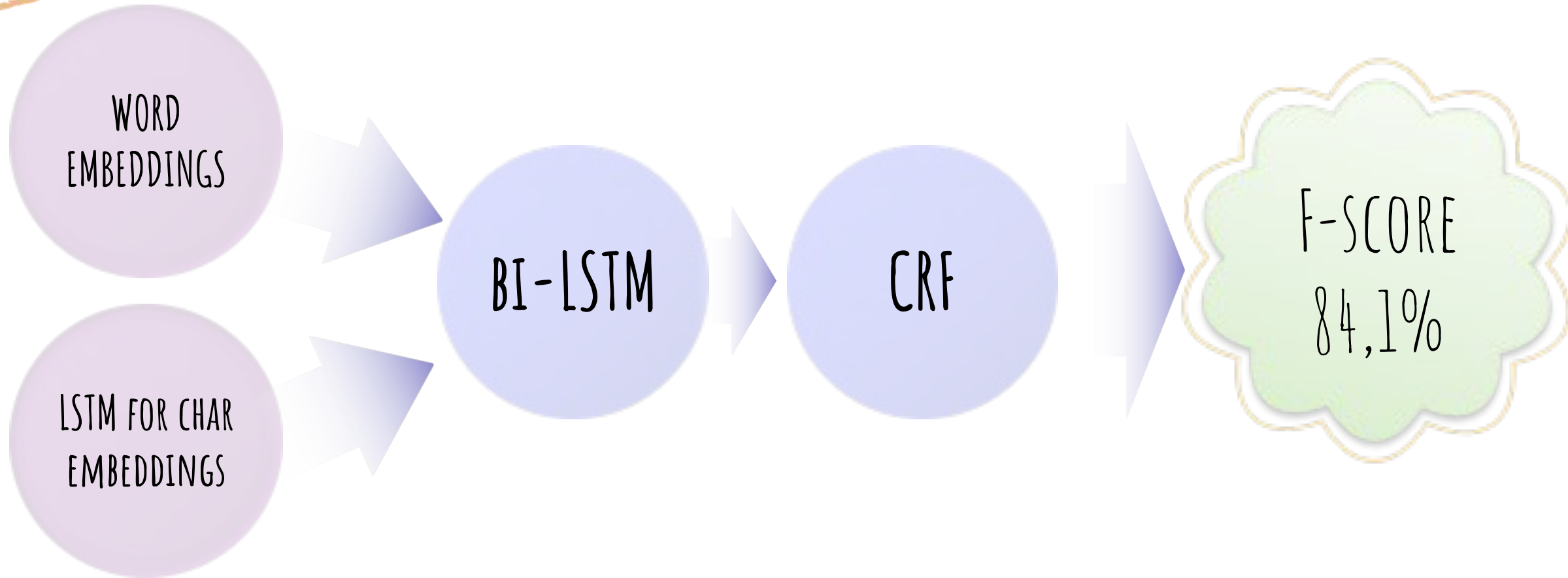
Time for one epoch: 85 seconds

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. RNN - LSTM
5. LSTM-CRF
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# LSTM-CNN-CRF

WORD EMBEDDINGS

CNN for char embeddings

BI-LSTM

CRF

F-score 83,9%

Time for one epoch: 130 seconds

# LSTM-LSTM-CRF

WORD EMBEDDINGS

LSTM FOR CHAR EMBEDDINGS

BI-LSTM

CRF

F-SCORE 84,1%

Time for one epoch: 305 seconds

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. RNN - LSTM
5. LSTM-CRF
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# BERT

| Fine-Tuning AdamW | | |
|---|---|---|
| #Words | Accuracy | F1 |
| 10000 | 0.87167 | 0.35335 |
| 100000 | 0.94249 | 0.7303 |
| All | 0.96268 | 0.83785 |

- Huggingface Library
- Basic Cased Model
- 4 Epochs
- AdamW Optimizer
  - Learning Rate: 3e-5
  - Epsilon 1e-8
- In the original Paper: F1 Score ~ 92%

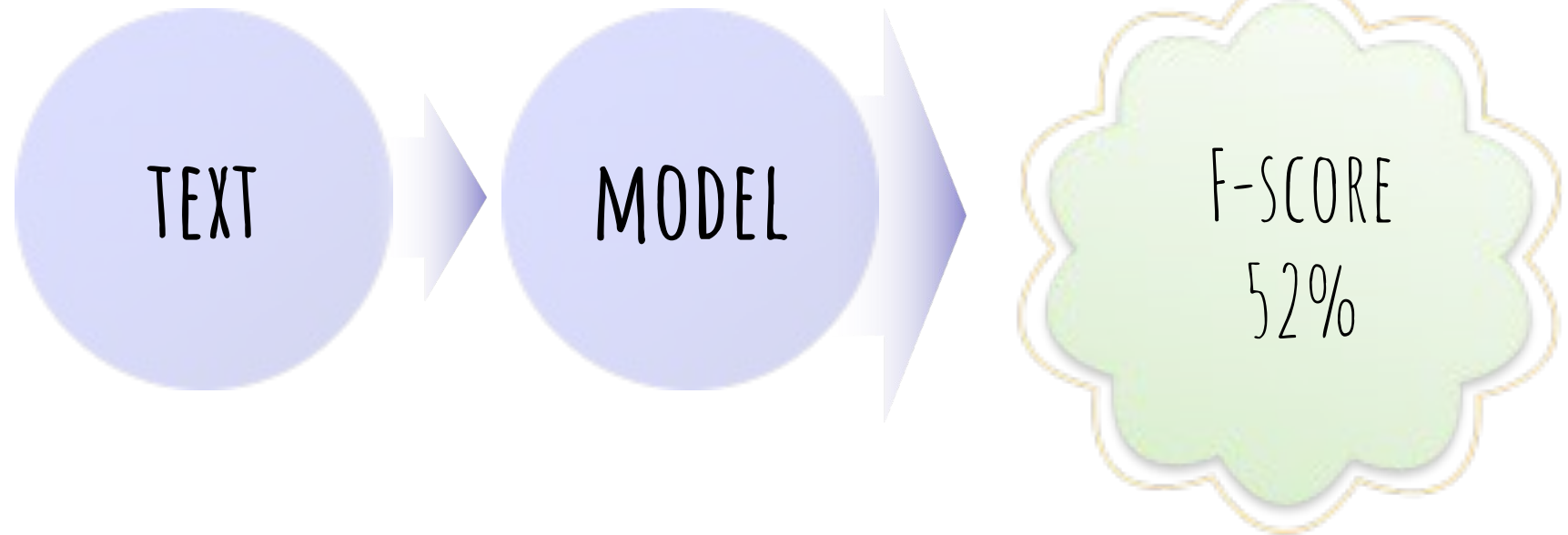| No Fine-Tuning AdamW | | |
|---|---|---|
| #Words | Accuracy | F1 |
| 10000 | 0.08524 | 0.0254 |
| 100000 | 0.80744 | 0.00514 |
| All | 0.84254 | 0.27026 |

Max
F-score
85.44%

# PRETRAINED BERT-NER

- Input -> Text
- Benchmark
- huggingface library
- baseline
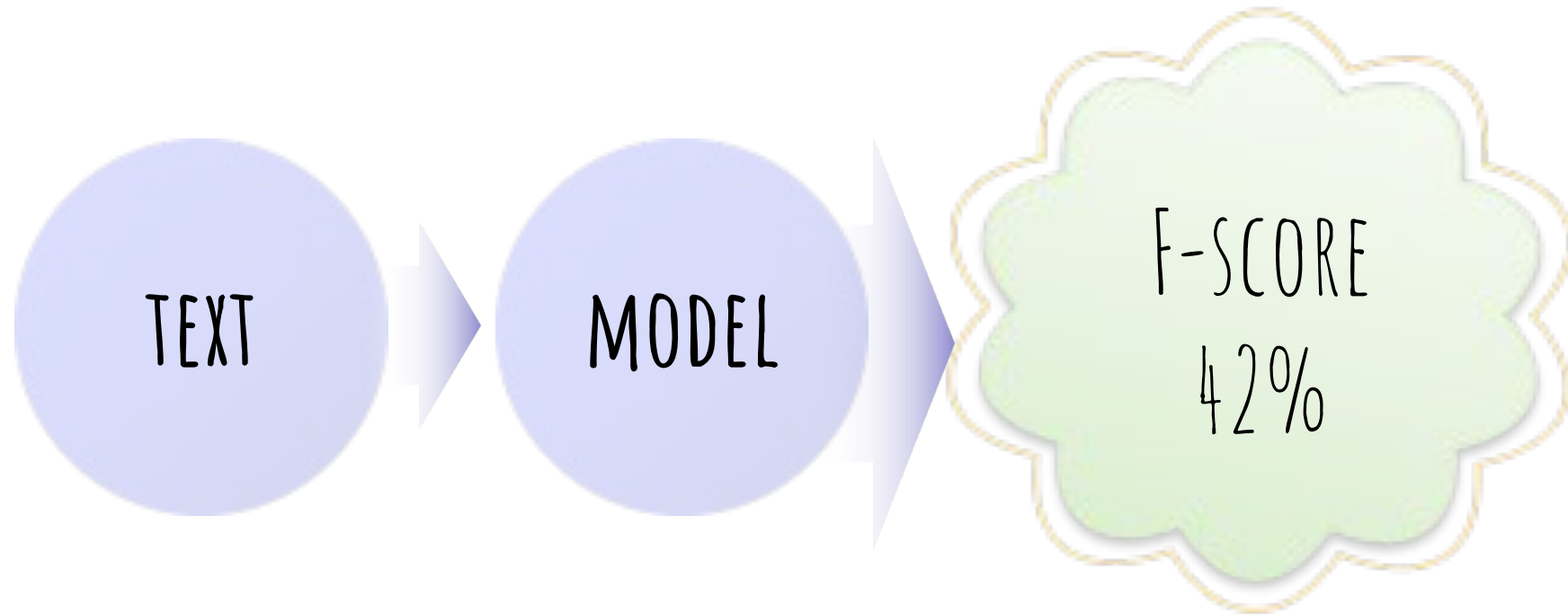
TEXT → MODEL → F-score 52%

# AGENDA

1. Motivation
2. Structure and size of the data set
3. CRF
4. RNN - LSTM
5. LSTM-CRF
6. Character embeddings
7. LSTM-CNN-CRF
8. BERT
9. SpaCy

# SPACY

- Input -> Text
- Benchmark
- NLP library
- Baseline

**TEXT** → **MODEL** → **F-SCORE 42%**

# Thank you for attention!

## Do you have any questions?