# Machine learning for natural language processing tasks

Aleksey Kulnevich,

Vladislav Radishevskii

13 September 2018
Dubna, Russia

# 1

## Introduction

# Machine Learning

**Machine learning** is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve result from experience without being explicitly programmed.

**Natural Language Processing**, or NLP, is the sub-field of AI that is focused on enabling computers to understand and process human languages.

London is the capital and most populous city of England and the United Kingdom.
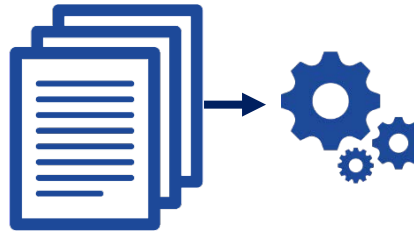
Geographic Entity         Geographic Entity         Geographic Entity

London is the capital and most populous city of England and the United Kingdom. Standing on the River Thames in the south east of the island of Great Britain, London has been a major settlement for two millennia. It was founded by the Romans, who named it Londinium.
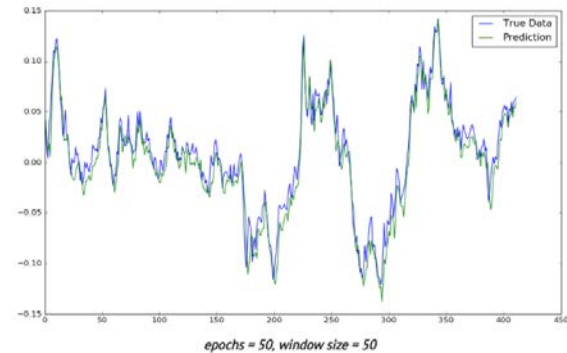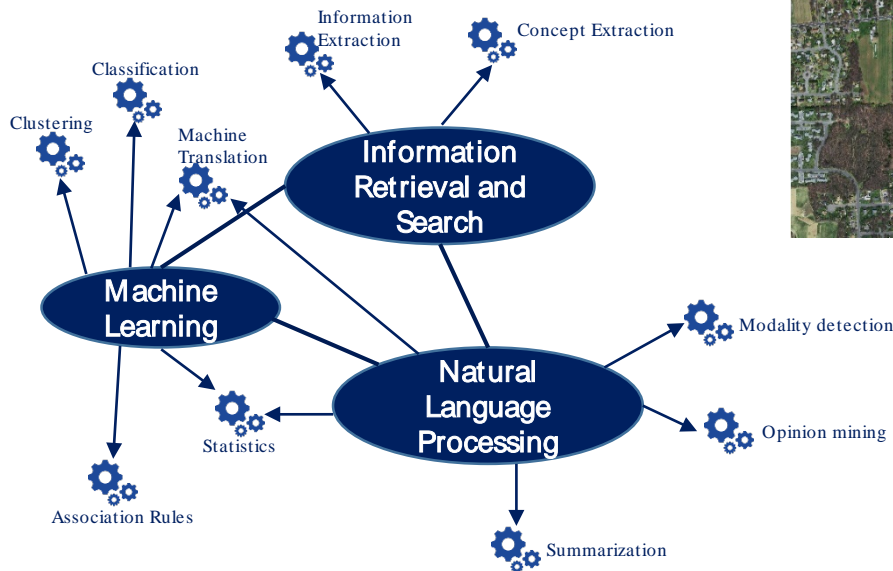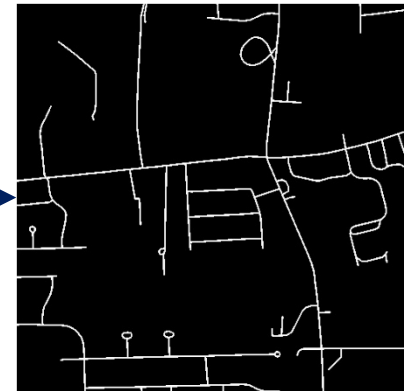
# Machine Learning

## ML TASKS
- ❖ *Information Extraction*
- ❖ *Machine Translation*
- ❖ *Image segmentation*
- ❖ *Object detection*
- ❖ *Gap filling*
- ❖ *Predictive analytics*
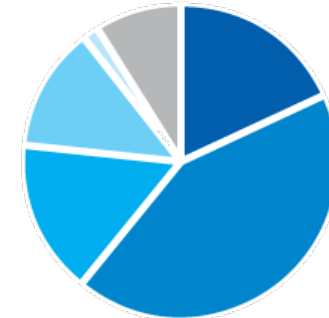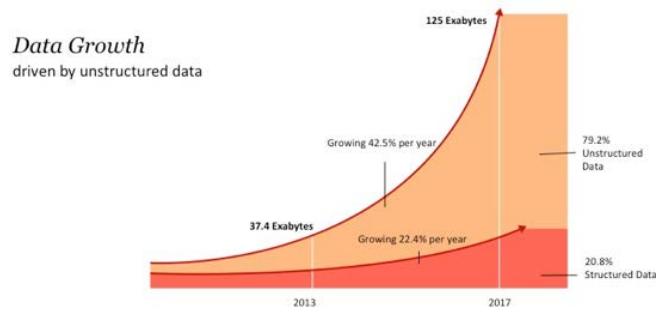


**Document structure:**
- ❖ *Named entity recognition*
- ❖ *Coreference resolution*
- ❖ *Type*
- ❖ *Keywords*
- ❖ *Attributes*





*epochs = 50, window size = 50*

# The importance of Text Analytics

**Structured / Unstructured Text Data**

❖ *Structured data represents only **20%** of the information available to an organization*
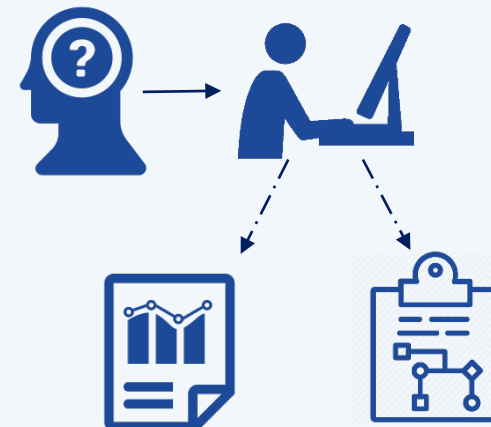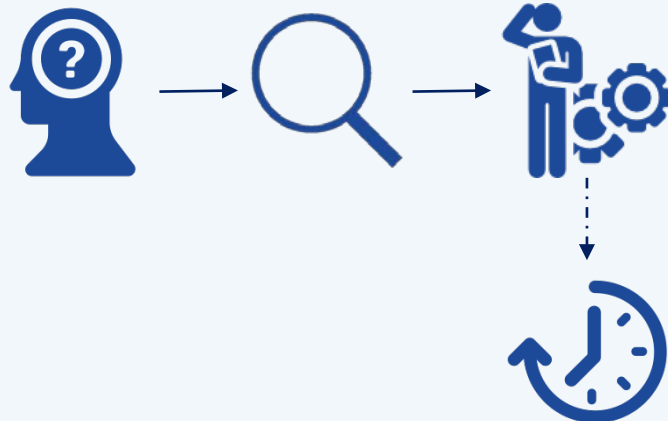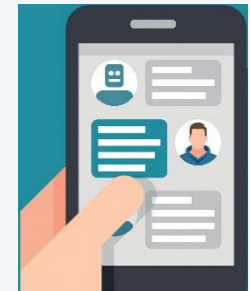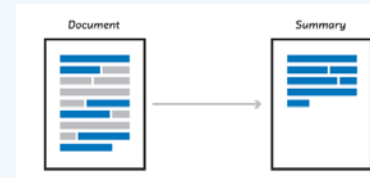❖ ***80%** of all the data is in unstructured form*



Data Growth
driven by unstructured data

125 Exabytes

Growing 42.5% per year

79.2% Unstructured Data

37.4 Exabytes

Growing 22.4% per year

20.8% Structured Data
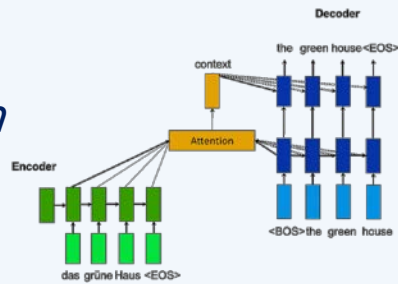
2013    2017

| | |
|---|---|
| Exclusively in structured format | 18.2% |
| Mostly in structured format | 42.9% |
| Equal split of structured and unstructured | 15.8% |
| Mostly in unstructured format | 12.8% |
| Exclusively in unstructured format | 1.6% |
| No clear understanding/Unsure | 8.7% |

❖ *If structured data is big, then unstructured data is huge*
❖ *Text analytics is the science of turning unstructured text into structured data*

# Natural Language Processing

## NLP TASKS

- ❖ *Named Entity Recognition*
- ❖ *Coreference Resolution*
- ❖ *Neural Machine Translation*
- ❖ *Chatbots*
- ❖ *Summary Extraction*
- ❖ *Answering Questions*
- ❖ *Ontology building*

# Objective of work

❖ *Feature extraction*
Word embeddings, char embeddings, morphological and additional tags
❖ *Building machine learning model for Named Entity Recognition*
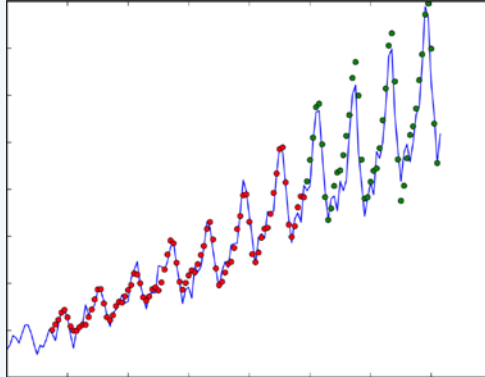Hybrid approach Bi-LSTM + CRF model
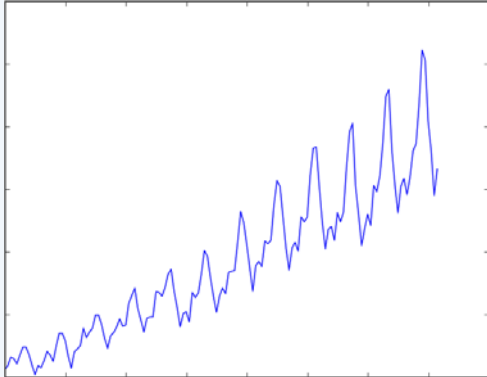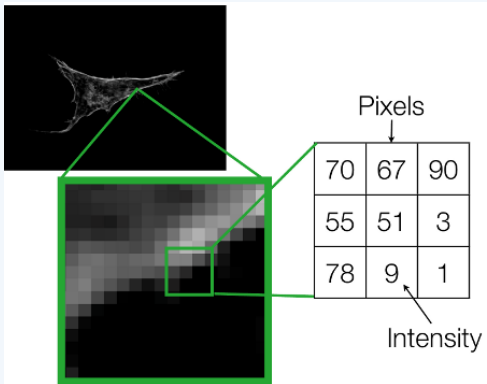❖ *Building machine learning model for Coreference Resolution*
Bi-LSTM model

2

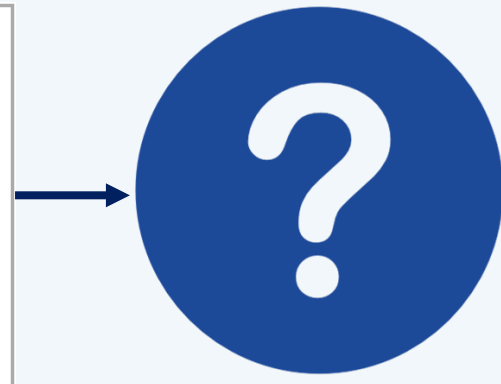# Vectorization problem

# Natural Language Processing

$$x = \{x_0, x_1.., x_n\}$$
$$y = x_{n+1}$$
$$x = \{x_{n+2}, x_{n+3}.., x_{n+k}\}$$
$$y = x_{n+k+1}$$
$$...$$

# Natural Language Processing

## Word2Vec – The Skip-Gram Model



### Intuition

*If two different words have very similar "contexts" (that is, what words are likely to appear around them), then the model needs to output very similar results for these two words.*

# 3

# Named Entity Recognition

# Natural Language Processing

**Entity** is concrete object of some type. For example, **Geoffrey Hinton** is an **entity** of *type "Person"*.

At the W party [Date] Thursday night [Time] at Chateau Marmont [Location], Cate Blanchett [Person] barely made it up in the elevator.

*At the W party* <Date> *Thursday* </Date> <Time>*night*</Time> *at* <Location> *Chateau Marmont* </Location>, <Person> *Cate Blanchett* </Person> *barely made it up in the elevator.*

"There was nothing about this storm that was as expected," said Jeff Masters, a meteorologist and founder of Weather Underground. "Irma could have been so much worse. If it had traveled 20 miles north of the coast of Cuba, you'd have been looking at a (Category) 5 instead of a (Category) 3."

Person    Organization    Location

# Natural Language Processing

**Named Entity Recognition**

**Named Entity Recognition is** subtask of information extraction that seeks to locate and classify *named entities* in text into pre-defined categories.
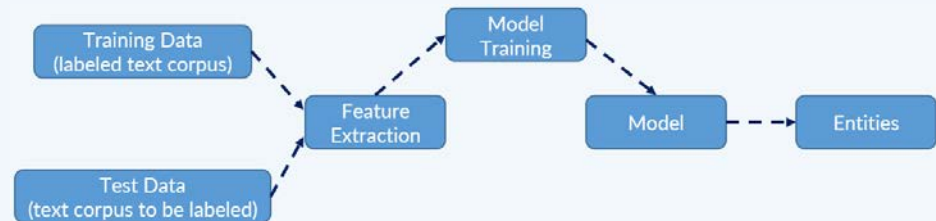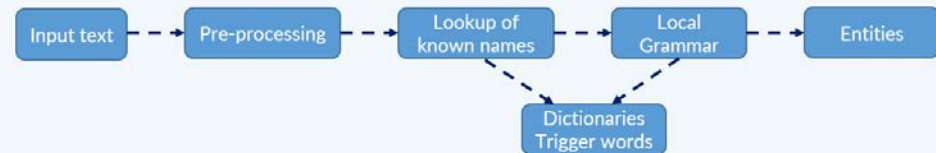
ML Approach
+ Flexibility
- Data for training

RB Approach
+ The ability to quickly find certain type of entities
- The need for specialized knowledge of linguistics
- Multiple rules

# Natural Language Processing

## Named Entity Recognition



$$p(\mathbf{y}|\mathbf{x}) = \frac{e^{Score(\mathbf{x},\mathbf{y})}}{\sum_{\mathbf{y}'} e^{Score(\mathbf{x},\mathbf{y}')}},$$

$$Score(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{T} A_{y_i,y_{i+1}} + \sum_{i=1}^{T} P_{i,y_i},$$

*Combination* of *CRF* model with a *Bi-LSTM* neural network encoding should increase the accuracy of the tagging decisions

*The CRF model is trained to predict a vector y = {$y_0$, $y_1$.., $y_T$} or tags given a sentence x = {$x_0$, $x_1$.., $x_T$}.*
*where A represents score of transition from tag I to tag j*
*P represents score of the $j^{th}$ tag of the word $i^{th}$*

# Natural Language Processing

**Named Entity Recognition**

**Features for improving accuracy:**
- ❖ Word embedding
- ❖ Char embedding
- ❖ Morphological features
- ❖ Additional tags: GEO, Orgn, Trad tm

**Language**: Python
**Frameworks**: NLTK, Numpy, Keras, Tensorflow
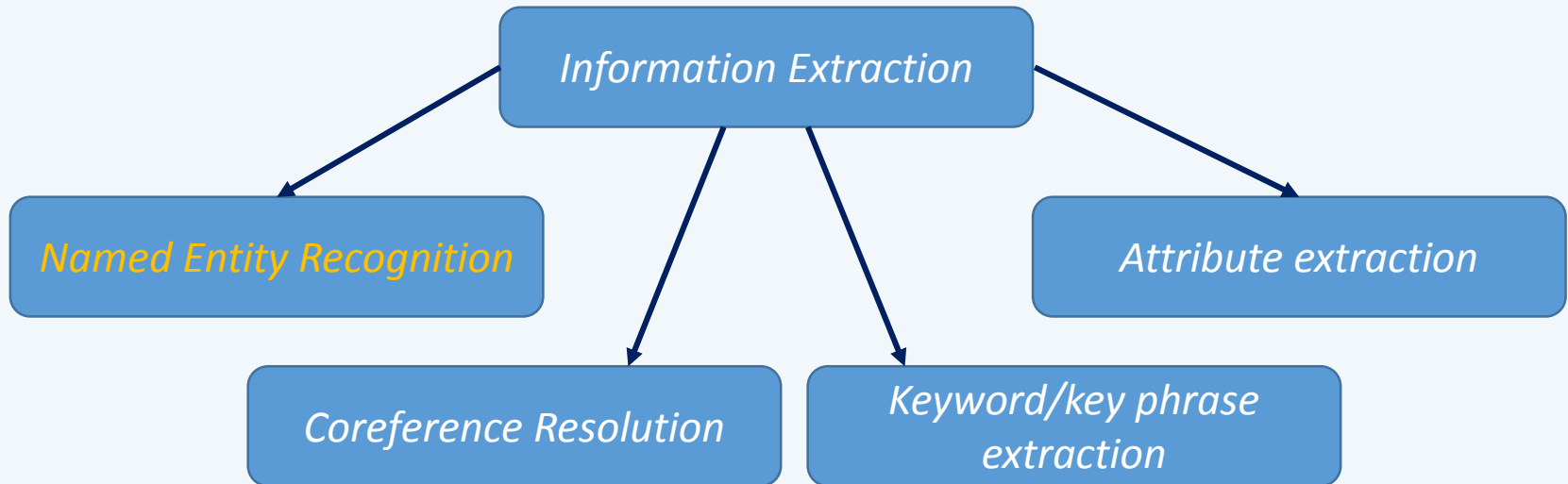
Natural Language Analysis with Python NLTK

K Keras

# Natural Language Processing

**Named Entity Recognition**

| Entity type | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| I-ORG | 0.81 | 0.80 | 0.81 | 1158 |
| B-PROD | 0.71 | 0.61 | 0.66 | 1590 |
| B-LOC | 0.82 | 0.85 | 0.83 | 1257 |
| I-LOC | 0.75 | 0.78 | 0.76 | 529 |
| I-PER | 0.89 | 0.87 | 0.88 | 919 |
| B-ORG | 0.78 | 0.73 | 0.76 | 1056 |
| I-PROD | 0.68 | 0.59 | 0.63 | 371 |
| B-PER | 0.85 | 0.86 | 0.86 | 711 |
| I-DATE | 0.97 | 0.98 | 0.97 | 955 |
| B-DATE | 0.91 | 0.92 | 0.91 | 749 |

# Natural Language Processing

**Information Extraction**



Information Extraction

Named Entity Recognition

Coreference Resolution

Keyword/key phrase extraction

Attribute extraction

- ❖ **Named entities** – *objects of specific types*
- ❖ **Coreference** – *chain of mentioning the named entity*
- ❖ **Keywords** – *are ideas and concepts that define what content is about*
- ❖ **Attributes** – *properties of objects*

# 4

## Coreference Resolution

# Coreference Resolution

**Coreference**, sometimes written **co-reference**, occurs when two or more expressions in a text refer to the same person or thing; they have the same referent.

FC Barcelona president Joan Laporta has warned Chelsea off star strike Lionel Messi.

This warning has generated dicouragement in Chelsea.

Aware of Chelsea owner Roman Abramovich's interest in the young Argentine, Laporta said last night: " I will answer as always, Messi is not for sale and we do not want to let him go."

**Coreference Resolution –** process of determining which mentions in a discourse refer to the same entity.

# Types of Coreference

## Anaphora

- ❖ **The music** was so loud that **it** couldn't be enjoyed.
- ❖ **Our neighbors** dislike the music. If **they** are angry, the cops will show up soon.

## Cataphora

- ❖ If **they** are angry about the music, **the neighbors** will call the cops.
- ❖ Despite **her** difficulty, **Wilma** came to understand the point.

## Split antecedents

- ❖ **Carol** told **Bob** to attend the party. **They** arrived together.
- ❖ When **Carol** helps **Bob** and **Bob** helps **Carol**, **they** can accomplish any task.

## Noun phrases

- ❖ Queen Elizabeth set about transforming **her husband**, **King George VI**, into a viable monarch. **Lionel Logue, a renowned speech therapist**, was summoned to help **the King** overcome his speech impediment.

# Supervised Approach
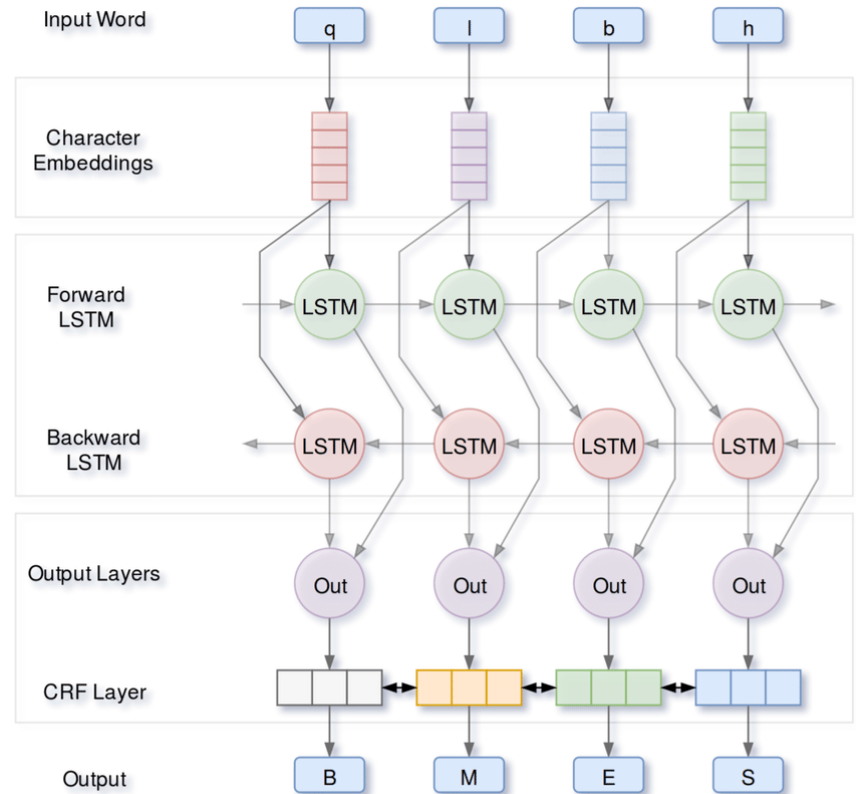
Based mainly on two methods:

❖ Binary classification

❖ Ranking method
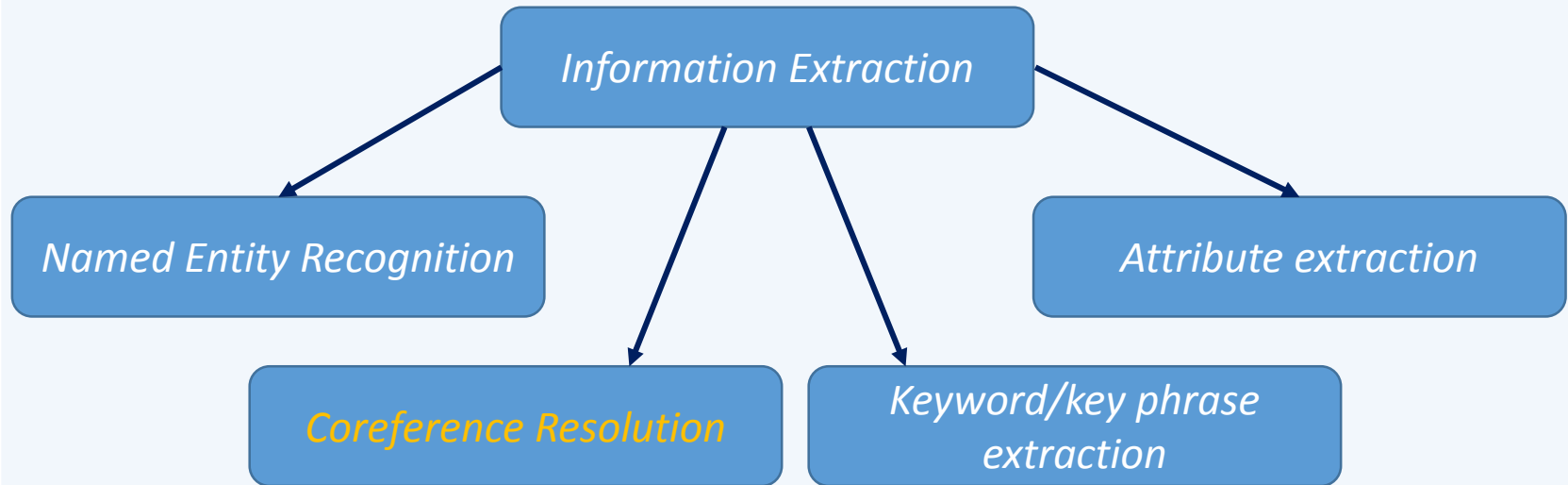
Pros:

❖ Learning algorithms usually generalize well

Cons:

❖ Quality is limited to quantity and quality of data

❖ Requires labeled data

# Natural Language Processing

```
Information Extraction

Named Entity Recognition          Attribute extraction

Coreference Resolution       Keyword/key phrase
                             extraction
```

- ❖ **Named entities** – *objects of specific types*
- ❖ **Coreference** – *chain of mentioning the named entity*
- ❖ **Keywords** – *are ideas and concepts that define what content is about*
- ❖ **Attributes** – *properties of objects*

econophysica

# Conclusion

We created neural network for solving named entity recognition and Coreference resolution problems:

❖ **Bi-LSTM Neural Network + CRF layer**

Input: word embeddings, char embeddings, morphological tags + geo tags + extra tags

❖ **Bi-LSTM Neural Network**

Input: word embeddings, char embeddings, morphological tags and names of entities (result of previous neural network)

econophysica

# Conclusion

This functionality is part of the text analytics system.
Visually it looks like this:

# Thanks for your attention!

Aleksey.Kulnevich@econophysica.com

a.d.kulnevich@gmail.com