

# Analiza koreferencji

(ang. coreference resolution)

# Koreferencja

- To zjawisko **współ odnoszenia się wyrażen referencyjnych** do tych samych **obiektów**.
- Wyrażenia referencyjne nazywane są **wzmiankami**.
- Grupa wzmianek odnosząca się do tego samego obiektu to **klaster koreferencyjny**.

„**Henryk Sienkiewicz** jest autorem ponad 20 nowel. (...) **Pisarz** jest także autorem powieści oraz reportaży.”

# Analiza koreferencji

- Celem **analizy koreferencji** jest połączenie **wyrażeń referencyjnych** (tzw. **wzmianek**) w grupy (tzw. **klastry koreferencyjne**) odnoszących się do tych samych obiektów.
- **Trudność zadania** wynika z faktu, że wzmianki należące do tego samego klastra **mogą mieć bardzo zróżnicowaną postać** (pod względem formy i kategorii). Wiąże się to m.in. ze względami stylistycznymi oraz praktycznymi.

„Tomek urodził się w Krakowie. Jak miał 16 lat, to jego rodzice zdecydowali się przeprowadzić do Szczecina. Długo nie mógł się pogodzić z ich decyzją o przeprowadzce. Tęsknił za miastem, w którym się urodził.”

# Typy wzmianek

- fraza rzeczownikowa
  - nazwy własne — Jan Nowak, Warszawa,
  - frazy pospolite — kierownik, stolica,
- zaimek
  - osobowy — on, ona, oni,
  - wskazujący — ten, ta, to, tamten,
  - dzierżawczy — jego, jej, ich,
- podmiot domyślny (Ø)
  - „Mężczyzna szedł po zmierzchu wzdłuż drogi.  
Przejeżdżający samochód oślepił go, przez co Ø **przewrócił** się do rowu.”

# Typologia relacji referencyjnych

- relacje bezpośrednie — wzmianka jawnie odnosi się do obiektu,
  - „Wczoraj zacząłem czytać Władcę Pierścieni.  
Ta książka bardzo mnie wciągnęła.”
- relacje pośrednie — wzmianka odnosi się do obiektu poprzez inny obiekt,
  - „Wczoraj zacząłem czytać Władcę Pierścieni.  
Przeczytałem już pierwszy rozdział.”

# Zastosowanie analizy koreferencji

- „Aparator wypłaci 1,20 zł dywidendy i skupi akcje własne za 20 mln zł.”
- „Aparator podjął decyzję o wypłacie dywidendy w drugiej połowie 2020 roku. Spółka wypłaci 1,20 zł za akcję.”

# Koreferencja w spaCy

# Pakiet **neuralcoref**

<https://github.com/huggingface/neuralcoref>

## ✨ **NeuralCoref 4.0: Coreference Resolution in spaCy with Neural Networks.**

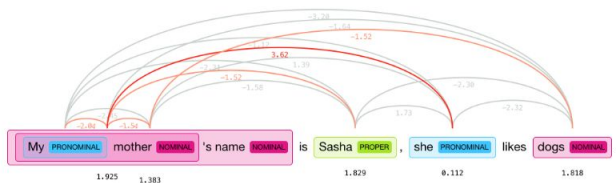
NeuralCoref is a pipeline extension for spaCy 2.1+ which annotates and resolves coreference clusters using a neural network. NeuralCoref is production-ready, integrated in spaCy's NLP pipeline and extensible to new training datasets.

For a brief introduction to coreference resolution and NeuralCoref, please refer to our [blog post](#). NeuralCoref is written in Python/Cython and comes with a pre-trained statistical model for **English only**.

NeuralCoref is accompanied by a visualization client [NeuralCoref-Viz](#), a web interface powered by a REST server that can be [tried online](#). NeuralCoref is released under the MIT license.

✨ Version 4.0 out now! Available on pip and compatible with SpaCy 2.1+.

release **v4.0.0** made with ♥ and **spaCy** build **passing**



```
# Load your usual SpaCy model (one of SpaCy English models)
import spacy
nlp = spacy.load('en')

# Add neural coref to SpaCy's pipe
import neuralcoref
neuralcoref.add_to_pipe(nlp)

# You're done. You can now use NeuralCoref as you usually manipulate a SpaCy document annotations.
doc = nlp(u'My sister has a dog. She loves him.')

doc._.has_coref
doc._.coref_clusters
```

Attribute	Type	Description
<code>doc._.has_coref</code>	boolean	Has any coreference has been resolved in the Doc
<code>doc._.coref_clusters</code>	list of Cluster	All the clusters of corefering mentions in the doc
<code>doc._.coref_resolved</code>	unicode	Unicode representation of the doc where each corefering mention is replaced by the main mention in the associated cluster.
<code>doc._.coref_scores</code>	Dict of Dict	Scores of the coreference resolution between mentions.
<code>span._.is_coref</code>	boolean	Whether the span has at least one corefering mention
<code>span._.coref_cluster</code>	Cluster	Cluster of mentions that corefer with the span
<code>span._.coref_scores</code>	Dict	Scores of the coreference resolution of & span with other mentions (if applicable).
<code>token._.in_coref</code>	boolean	Whether the token is inside at least one corefering mention
<code>token._.coref_clusters</code>	list of Cluster	All the clusters of corefering mentions that contains the token

- **Operating system:** macOS / OS X · Linux · Windows (Cygwin, MinGW, Visual Studio)
- **Python version:** Python 3.6+ (only 64 bit)
- **Package managers:** [pip]



# Pakiet **neuralcoref** — parametry

Parameter	Type	Description
<code>greedyness</code>	float	A number between 0 and 1 determining how greedy the model is about making coreference decisions (more greedy means more coreference links). The default value is 0.5.
<code>max_dist</code>	int	How many mentions back to look when considering possible antecedents of the current mention. Decreasing the value will cause the system to run faster but less accurately. The default value is 50.
<code>max_dist_match</code>	int	The system will consider linking the current mention to a preceding one further than <code>max_dist</code> away if they share a noun or proper noun. In this case, it looks <code>max_dist_match</code> away instead. The default value is 500.
<code>blacklist</code>	boolean	Should the system resolve coreferences for pronouns in the following list: ["i", "me", "my", "you", "your"]. The default value is True (coreference resolved).
<code>store_scores</code>	boolean	Should the system store the scores for the coreferences in annotations. The default value is True.
<code>conv_dict</code>	dict(str, list(str))	A conversion dictionary that you can use to replace the embeddings of <i>rare words</i> (keys) by an average of the embeddings of a list of <i>common words</i> (values). Ex: <code>conv_dict={"Angela": ["woman", "girl"]}</code> will help resolving coreferences for <i>Angela</i> by using the embeddings for the more common <i>woman</i> and <i>girl</i> instead of the embedding of <i>Angela</i> . This currently only works for single words (not for words groups).

```
import spacy
import neuralcoref

nlp = spacy.load('en')

# Let's try before using the conversion dictionary:
neuralcoref.add_to_pipe(nlp)
doc = nlp(u'Deepika has a dog. She loves him. The movie star has always been fond of animals')
doc._.coref_clusters
doc._.coref_resolved
# >>> [Deepika: [Deepika, She, him, The movie star]]
# >>> 'Deepika has a dog. Deepika loves Deepika. Deepika has always been fond of animals'
# >>> Not very good...

# Here are three ways we can add the conversion dictionary
nlp.remove_pipe("neuralcoref")
neuralcoref.add_to_pipe(nlp, conv_dict={'Deepika': ['woman', 'actress']})
# or
nlp.remove_pipe("neuralcoref")
coref = neuralcoref.NeuralCoref(nlp.vocab, conv_dict={'Deepika': ['woman', 'actress']})
nlp.add_pipe(coref, name='neuralcoref')
# or after NeuralCoref is already in SpaCy's pipe, by modifying NeuralCoref in the pipeline
nlp.get_pipe('neuralcoref').set_conv_dict({'Deepika': ['woman', 'actress']})

# Let's try again with the conversion dictionary:
doc = nlp(u'Deepika has a dog. She loves him. The movie star has always been fond of animals')
doc._.coref_clusters
# >>> [Deepika: [Deepika, She, The movie star], a dog: [a dog, him]]
# >>> 'Deepika has a dog. Deepika loves a dog. Deepika has always been fond of animals'
# >>> A lot better!
```

# Koreferencja dla j. polskiego

# Narzędzia dla j. polskiego

<http://zil.ipipan.waw.pl/PolishCoreferenceTools>



PolishCoreferenceTools

## Menu

Ling. Engineering Group  
IPI PAN  
CLIP  
Facebook  
YouTube channel

## Polish Coreference Tools

This page describes the tools created as part of the [CORE](#) project.

### Mention detection

- [MentionDetector](#) – mention detection tool

### Coreference resolution

- [Ruler](#) – a rule-based coreference resolution tool
- [Bartek](#) – a statistical coreference resolution tool

### Evaluation

- [Scoreference](#) – a mention detection and coreference resolution evaluation tool

### Coreference annotation

- [DistSys](#) – a distribution system for texts for any kind of manual annotation
- [MMAX4CORE](#) – a modified version of the MMAX2 annotation tools, adjusted for the needs of the [CORE](#) project

### Converters

- [PCC converters](#) – a suite of converters between data formats used in [CORE](#) project

### Visualization

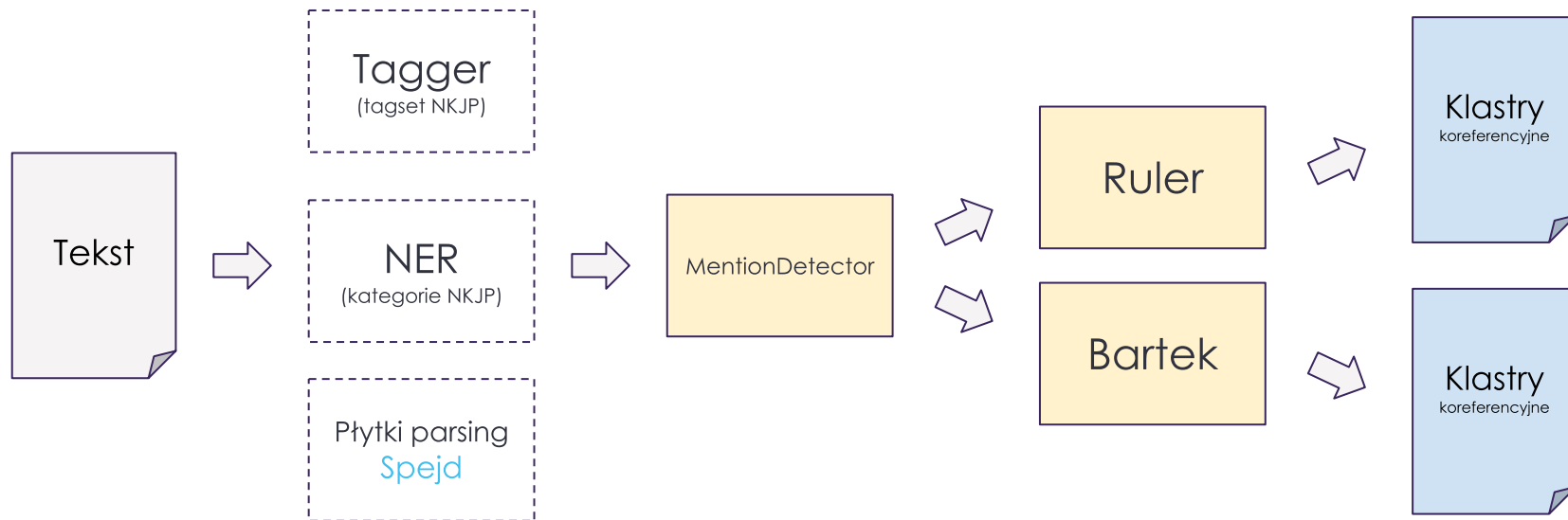
- [brat4core](#) – a modified version of [brat annotation environment](#), tailored for visualization of coreference in PCC.

### Other

- [Text Selector](#) – tool for manual text inspection and selection

- wzmianki
  - [MentionDetector](#)
- klastry koreferencyjne
  - [Ruler](#) — regułowcy
  - [Bartek](#) — statystyczny

# Potok przetwarzania



# Koreferencja przez Multiservice

<http://multiservice.nlp.ipipan.waw.pl>

Multiservice Demo

Uruchom przetwarzanie

Wynik ostatniego ządania

Zgłoś problem

Pomoc

Segmentacja

Morfologia/składnia

Słowa

Grupy

Jednostki nazwane

Wzmianki

Koreferencje

1 To będzie już druga próba licytacji nieruchomości na pl . Słonecznym , którą urzędnicy wytropili po latach poszukiwań majątku Adama Gesslera .

2 Jego dług wobec miasta szacują dziś na ok . 27 mln zł . Już w 1992 r . , wkrótce po podpisaniu umowy najmu lokalu na Rynku Staromiejskim , zaczęły się problemy z czynszem . Sąd orzekł eksmisję . Dotąd miastu udało się odzyskać ledwie kilkadziesiąt tysięcy złotych długu .

3 Sprawa budzi wielkie emocje , bo choć Adam Gessler jest słynnym restauratorem , oficjalnie nie ma nic . Nawet wynajęta przez Zakład Gospodarowania Nieruchomościami w Śródmieściu firma detektywistyczna nie znalazła majątku .

4 Pozostają dwa mieszkania na Żoliborzu , wyceniane przed rokiem na blisko 4.3 mln zł . Będą licytowane za dwie trzecie ceny . W ZGN wymislił , żeby miasto przystąpiło do licytacji . Jeśli uda się kupić nieruchomość , komornik pospłaca wierzycieli Adama i Piotra Gesslerów . A miasto będzie mogło w przyszłości sprzedać korzystnie atrakcyjny dom .

5 Licytacje odbędą się w środę . - Korzyści z wycycytowania domu będą niewielkie w stosunku do ogromnego długu pana Gesslera . Chodzi jednak o to , żeby wiedział , że miasto nie zrezygnuje z upominania się o swoje - tłumaczyła " Gazecie " Małgorzata Mazur , dyrektorka ZGN .

# Literatura

Maciej Ogrodniczuk (2019). **Automatyczne wykrywanie nominalnych zależności referencyjnych w polskich tekstach współczesnych.**

Wydawnictwo Uniwersytetu Warszawskiego

[https://www.wuw.pl/data/include/cms/Automatyczne\\_wykrywanie\\_nominalnych\\_Ogrodniczuk\\_Maciej\\_2019.pdf](https://www.wuw.pl/data/include/cms/Automatyczne_wykrywanie_nominalnych_Ogrodniczuk_Maciej_2019.pdf)