

Inforex — a Collaborative System for Text Corpora Annotation and Analysis

Michał Marcińczuk

Marcin Oleksy

Jan Kocoń

G4.19 Research Group

Department of Computational Intelligence

Faculty of Computer Science and Management

Wrocław University of Technology, Wrocław, Poland

{michal.marcinczuk, marcin.oleksy, jan.kocon}@pwr.edu.pl

Abstract

We report a first major upgrade of Inforex — a web-based system for qualitative and collaborative text corpora annotation and analysis. Inforex is a part of Polish CLARIN infrastructure¹. It is integrated with a digital repository for storing and publishing language resources² and it allows to visualize, browse and annotate text corpora stored in the repository. As a result of a series of workshops for researchers in Humanities and Social Sciences we improved the graphical interface to make the system more friendly and readable for non-experienced users. We also implemented a new functionality for a gold standard annotation which includes private annotations and annotation agreement by a super-annotator.

1 Introduction

Digital humanities (DH) create new demand and challenges for development of new or existing tools and systems for text documents manipulation, processing, analysis and visualization. CLARIN-PL — the Polish part of CLARIN infrastructure — tries to rise the challenges associated with DH for Polish language. Among many other issues, there is a need for an intuitive and easy to use system for qualitative text corpora management, annotation, analysis and visualization. To fulfill these needs we develop such a system called Inforex. In this article we present the current state of the system development.

The decision to create a system for text corpora annotation was taken in 2009 when there were no such systems which support collaborative

work. On that time the only existing tools were desktop applications for individual work such as GATE (Cunningham et al., 2011) or Manufakturyzacja Luna (Marciniak et al., 2010). Since 2010 several systems have emerged, like WebAnno 3 (Eckart de Castilho et al., 2016) or GATE Teamware (Bontcheva et al., 2013).

The first version of Inforex system was released in 2010 and its initial role was to construct corpus-based linguistic resource for various tasks from the field of natural language processing, including named entity recognition (Marcińczuk et al., 2011), shallow parsing (Radziszewski and Piasecki, 2010), word sense disambiguation (Bas et al., 2008), recognition of semantic relations between named entities (Marcińczuk and Ptak, 2012). It was used to develop two major (at that time) resources for Polish: Corpus of Wrocław University of Technology called KWPr (Broda et al., 2012) (within the NEKST³ project) and Corpus of Economic News (CEN) (Marcińczuk et al., 2013) (within the SyNaT project⁴). Later, in 2013 Inforex was used to construct another major resource, which is Polish Corpus of Suicide Notes (PCSN)⁵ (Marcińczuk et al., 2011) guided by Monika Zaśko-Zielińska (2013). Until now the system has been used to access the corpus. The access is granted on a demand after obtaining a permission form Wrocław University.

In 2013 Poland joined CLARIN — European Research Infrastructure for Language Resources and Technology. The goal of CLARIN is to make the language technologies more accessible to researches from humanities and social sciences, which in most cases do not have the technical skills to use many of the tools on their own. At that time we made a decision to make Inforex a part

¹<http://clarin-pl.eu>

²<http://clarin-pl.eu/dspace>

³<http://nekst.ipipan.waw.pl/>

⁴<http://www.synat.pl/>

⁵<http://pcsn.uni.wroc.pl/>

of the Polish CLARIN infrastructure. In 2015–2017 we have organized several workshops for researchers in humanities and social sciences. The workshops showed us several user experience issues. System GUI turned out to be not enough intuitive for non-experienced users. Then, first of all, it needed to be simplified. Second problem was connected with the methodology. The researchers use various tools for corpora analysis (including spreadsheets) and Inforex may be treated as some kind of pre-processing tool that allows to prepare corpus for further analysis. Data export was possible but complicated and required an access to a database. Users feedback proved that the easy form of data export is one of the crucial needs. After the set of workshops we gathered more information about other important needs (also in the form of questionnaires) like access to a custom annotation schemas definition or data visualisation. Some of them have been already implemented and the other are under construction.

2 Inforex Features Overview

In the following sections we present the main functionalities and features of the Inforex system.

2.1 Web-based Access

Inforex is a web-based tool which does not require installation. It can be accessed by any web-browser which support JavaScript. Despite Inforex is built on several universal JavaScript libraries and frameworks (jQuery, jQuery extensions and Bootstrap) we suggest using Chrome and Firefox. These two web browsers are used to test the system on daily bases. Users might use other browsers as well, however we are not able to validate all functions in each of the available web browsers, thus some minor issues might occur.

2.2 Authorized and Public Access

Corpora stored in Inforex can be accessed by authorized and unauthorized users. The manager of the corpus (the owner or a user with specific privileges) decides what type of information from the corpora can be publicly available. For instance, only authorized users can have access to documents' content and can modify the corpus annotations while unauthorized users may have access to some statistics or annotation frequency lists.

2.3 Integration with DSpace as a Part of Polish CLARIN Infrastructure

Inforex system is available at <http://inforex.clarin-pl.eu> and it is part of Polish CLARIN infrastructure. This installation is integrated with the official repository for language resources in Polish CLARIN⁶. The repository runs on DSpace system⁷. When a user registers in <https://clarin-pl.eu/dspace/>, he also gains access to Inforex system. At this stage accounts are automatically synchronized. In the future both systems will use unified federation authorization.

2.4 Collaboration

Inforex offers several ways for collaborative work on a single corpus. One of them is the access to the same corpora for different authorized users. The other one is a selective, task-oriented access to the same document. For instance, different groups of users can have access to document's metadata. The last one is the "2+1" annotation, i.e. two or more users annotate the same set of documents independently and the super-annotator creates the final set of annotations based on their input. More about this type of collaboration is presented in Section 3.2.

2.5 Qualitative Document Annotation

Inforex was designed for qualitative document annotation. This means it does not offer a fast and robust search functions over large corpora with millions of documents. Such functionality can be obtained using other existing tools designed for it, for instance Sketch Engine (Kilgarriff et al., 2014) or NoSketch Engine (Rychlý, 2007). Inforex is suited for medium size corpora (containing thousands of small documents) and to manually describe documents in terms of their metadata, annotations (types of phrases organized in a hierarchy), annotation attributes, relations between annotations and annotation frames.

2.6 Language-independent

Inforex is language-independent in the sense that it can handle documents in any natural language. So far it has been used to annotate Polish, English and Hebrew texts (see Section 3.2).

⁶<https://clarin-pl.eu/dspace/>

⁷<https://github.com/ufal/clarin-dspace>

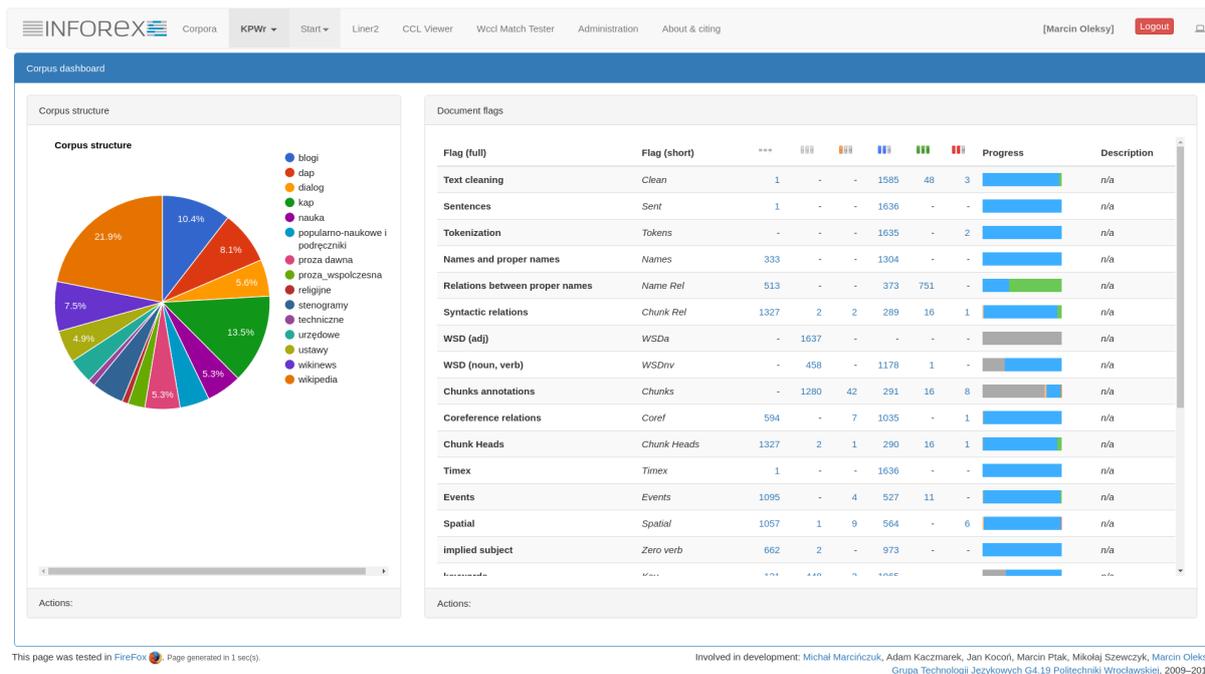


Figure 1: Corpus overview

2.7 Document Visualisation

Inforex can handle documents in two formats: plain text and XML. For XML documents it is possible to display their content in a visually formatted way. This allows to highlight the document structure what improves the user experience while browsing and annotating documents. Sample visualizations of different types of documents are presented in Figure 3.

2.8 Document Description

Inforex supports four types of information units which can be used to describe documents content:

1. Metadata — an information unit which is assigned to whole document (author name, document creation time, source, etc.).
2. Annotation — an information unit which is assigned to a sequence of words in the document content. Each annotation is described with a category (categories can be organized in a hierarchy) and a set of attributes. The set of attributes depends on the semantic interpretation of the annotation category. For instance, for named entities it can be a lemma, for temporal expressions it can be a normalized value of the expression and for event mentions it can be an event modality.

3. Relation — an information unit which is assigned to a pair of annotations. It is a directed link between two annotations of some category.
4. Frame — an information unit which is assigned to a set of annotations. Frame consists of a set of annotations with roles assigned to them. This type of structure can be used for event annotations (LCD, 2005).

3 Recent Improvements

In the following sections we present the recent major improvements of Inforex system.

3.1 Modern Layout

A set of workshops carried out from 2015 to 2017 showed that there was the need for an adjustment of user interface to a new group of users — researchers in humanities and social sciences not involved in NLP tools development. New users reported confusion with the large amount of information and the number of available functions. The need of interface simplification appeared while functionalities of the system would remain unchanged. Thus, Inforex layout has been upgraded and modernized. It involved not only a design lifting of the user interface but also changes in navigation panels. The comparison of *old* and *new*

The screenshot shows the Inforex web application interface. At the top, there is a navigation bar with the Inforex logo and various menu items like 'Corpora', 'KPWr', 'Documents', 'Liner2', 'CCL Viewer', 'Wccl Match Tester', 'Administration', 'About & citing', and a user profile 'Marcin Oleksy' with a 'Logout' button. Below this is a breadcrumb trail: '1 z 1637: wikipedia » Toronto Dominion Centre'. The main content area is divided into two panels. The left panel, titled 'Document content', shows a text snippet from Wikipedia about the Toronto Dominion Centre, with several words and phrases highlighted in yellow and blue, indicating annotations. The right panel, titled 'Annotation details', shows the details for a specific annotation: 'Id: 31822', 'Text: Toronto Dominion Centre', 'Type: nam_fac_goe', and 'Lemma: Toronto Dominion Centre'. Below this is a table for 'Annotation relations' with columns for 'Id', 'Relation type', and 'Target annotation'. At the bottom of the page, there is a footer with the text: 'This page was tested in FireFox... Page generated in 0 sec(s). Involved in development: Michał Marcińczuk, Adam Kack, Jan Kocoń, Marcin Ptak, Mikołaj Szewczyk, Marcin Oleksy Grupa Teorii Językowych G4.19 Politechniki Wrocławskiej, 2009–2017 Grupa Technologii Językowych G4.19 Politechniki Wrocławskiej, 2009–2017'.

Figure 2: Document annotation view

layout is presented in Figure 4.

3.2 Annotation Agreement

Reliability is a key value in the creation of a good quality corpora for learning and testing of NLP tools. The current version of Inforex enables simultaneous and independent annotation of the same text sample by more than one annotator. Moreover, the annotation process coordinator may keep track of inter-annotator agreement between two raters thanks to the *Agreement module* which uses Positive Specific Agreement (PSA) measure (Hripcsak and Rothschild, 2005) to calculate the reliability (see Figure 5). View configuration gives the opportunity to define annotation layers, subsets or categories, users and set of documents that have to be analysed. The coordinator may also specify a comparison mode: whether the system has to take into consideration the annotation boundaries only or boundaries and categories. It may also include annotation lemmas. Inter-annotator agreement is a very important indicator of the annotation guidelines clearness or cohesion. Keeping track of changes of the inter-annotator agreement between subsequent annotation iterations helps to improve the quality of the annotation guidelines. Agreement module makes that process easier and faster.

Inforex system also supports the curation of

the annotation process (see Figure 6). The curator can make choice between two different annotators choices, or even reject consistent but incorrect annotations. Thanks to that module several Gold Standard projects were performed e.g. Polish Coreference Corpus (Ogrodniczuk et al., 2015) for definite descriptions annotation and Polish Spatial Texts corpus for the annotation of dynamic spatial expressions.

4 Applications

In the following sections we present several practical applications of the Inforex system.

4.1 KPWr

KPWr (Polish Corpus of Wrocław University of Technology) (Broda et al., 2012) is a corpus of written and spoken documents available on the Creative Commons license which is intended primarily as a training and testing material for NLP tools being developed at Wrocław University of Science and Technology. It is successively enriched with annotation layers. Inforex recently supported manual text annotation within such layers as temporal expressions and their normalizations, events (and description of event attributes), spatial expressions and semantic roles. In order to prepare temporal expressions annotation (Kocoń et al., 2015) a new annotation scheme based on

Ewa Kaczmarska	2014-05-2 20:31
Ewa Kaczmarska	Hejka.
Dawid Maciejewski	2014-05-2 20:45
Dawid Maciejewski	Hej.
Ewa Kaczmarska	2014-05-2 20:50
Ewa Kaczmarska	Składaj, jako lipa. Sprawdzaliśmy dużo programów i nie padła, więc... my chyba i tak to zrobimy, a skłapy ma być do 20.

(a) Facebook conversation.

Wieża szybowa (górnictwo)

Wieża szybowa - wysoka konstrukcja żelbetowa, stalowa, betonowa, i zabudowane są koła kierownicze lin wyciągu szybowego lub maszyna głębenia szybu stosowane są wieże tymczasowe.

Rodzaje wież szybowych

Basztowe (wolnostojące)

- Trzonowe pełne
- Trzonowe dzielone
- Trzonowe słupowe
- Słupowe

Zastrzałowe

- Jednozastrzałowe
- Dwuzastrzałowe (koziłowe).

(b) Wikipedia article.

נפט' יו' א' טבת שנת ו'ור
יהודא מאת
בני אחיו לפק
פה נטמן איש הישר עושה טוב
רמרע סר : בחיותו על אדמתו
אכל מיגיע כפו : ועמד בצדקת
עד סופו : ה"ה האלוף הקצין כהר"ר
יהודא ליב ב"ה תנחום מיינשטר
זצל פ"ל : כל ימיו עסק באמונה
כל מפעלו : ולעניים ודלי קרובי
ורחוקי ניזון מלחמו : תמיד החזיק
מלמדו : ללמוד תורת ה' לבנים
שינת עדנה יישן גופו : בגן עדן
ישמח נשמתו : ושבק לן חייל : תחת
שכל מעשיו היו לש"ש : ודבק
את עצמו באלהים חיים : לפק
ת"ל צ"ב"ה

(c) Hebrew document.

Figure 3: Sample documents visualizations

TimeML was added. These categories refer to a date, time of a day, duration and frequency of an event. Annotation lemmas perspective was used to provide normalized temporal expressions, revealing that the term 'lemma' in Inforex may function as a broad concept. The Annotator perspective from the system also supports event annotation (Marcinićzuk et al., 2015). There are seven coarse-grained categories of events, i.e. action, state, reporting, perception, aspectual, intensional action and intensional state. The categorization was based on the TimeML guidelines with some modifications. It also involved creation of a new annotation scheme. The flexibility in adding new annotation layers (setting the new annotation categories) is one of the most important features. The possibility of establishing relations between annotated fragments is not less relevant. It was crucial e.g. for spatial expressions annotation. Its main goal was to extract different ways of distributing spatial information throughout a sentence by reviewing the lexical and grammatical signals of various relations between objects (Marcinićzuk et al., 2016).

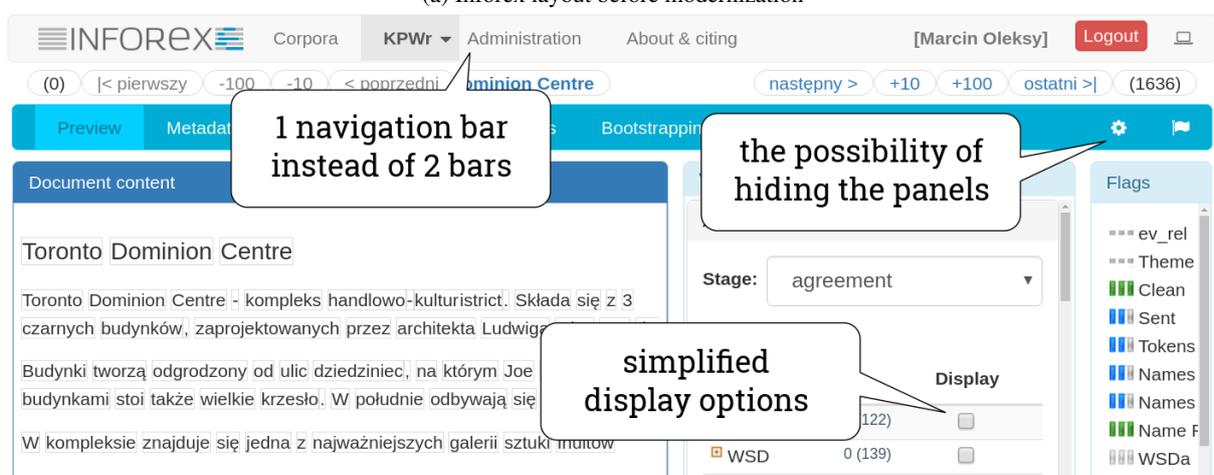
4.2 European Legal Texts

As practice shows, although Inforex was primarily developed for Polish language, that it can also be used to work with documents written in other languages. Inforex features and functionalities are useful e.g. in examining current EU official literature related to territorial development and urban planning. Authors of this analysis first uploaded EU Territorial Policy Documents 2007-2016⁸ to CLARIN-PL DSpace repository and then imported it to the Inforex system. The corpus was divided into 4 subcorpora and prepared for qualitative and quantitative analysis. The review of the key strands enabled the identification of its 8 core values (or principles) for further statistical and contextual analysis. After ascribing to each category its textual triggers (word forms), a quantitative analysis using words frequency lists generated by Inforex was performed. Manual annotation with a newly defined set of annotations and Annotation Browser with the possibility of exporting data were a great support for qualitative analysis — detailed contextual analysis of the corpus focused on two crucial categories: *Participation* and *Communication*.

⁸<http://hdl.handle.net/11321/316>



(a) Inforex layout before modernization



(b) Inforex layout after modernization

Figure 4: Inforex layouts comparison

4.3 Hebrew Corpus

Inforex supports manual annotation even if the text is written using non-latin alphabet and a right-to-left notation. One of the system applications was related to a corpus of Hebrew gravestone inscriptions. It also involved the creation of a new annotation schema. Categories referred mainly to the pragmatic level of communication (e.g. initial and final expressions, laudations, death circumstances). The perspective of annotation lemmas was used to enter Polish translations of annotated fragments, which also showed that the lemma attribute may be a broad term especially in the case of practical applications of the system.

4.4 Other Corpora

Inforex was used to prepare the training data during participation in BSNLP 2017 shared task on

multilingual named entity recognition aimed at recognizing mentions of named entities in web documents in Slavic languages, their normalization / lemmatization, and cross-language matching (Marciniuk et al., 2017). The system also supported the annotation of the corpora constructed specially for specific tasks from the field of natural language processing e.g. Polish Coreference Corpus for definite descriptions annotation and Polish Spatial Texts corpus for the annotation of dynamic spatial expressions. It involved creation of dedicated annotation layers but, what is important, in these tasks the new module of the system (Annotation Agreement and "2+1" annotation) was used for the first time, which significantly improved the time of preparation of annotated training and testing corpora.

5 Summary

Inforex system, as a part of CLARIN-PL infrastructure, is gradually developed. Although its initial role was to construct qualitative linguistic resources for various tasks from the field of natural language processing, recently it is also used by scientists for other purposes. We received an important and constructive feedback from users during and after workshops related to CLARIN-PL tools and resources. As users have different needs, we identified the common functionalities and implement them as soon as possible in order to boost their research tasks and provide new possibilities. We also challenged with the fact that many researches from the field of digital humanities are not experienced users of such systems and we made Inforex as easy and intuitive as possible.

Acknowledgments

Work financed as part of the investment in the CLARIN-PL research infrastructure funded by the Polish Ministry of Science and Higher Education.

References

- Dominik Bas, Bartosz Broda, and Maciej Piasecki. 2008. *Towards Word Sense Disambiguation of Polish*. In *Proceedings of the International Multiconference on Computer Science and Information Technology, {IMCSIT} 2008, Wisla, Poland, 20-22 October 2008*. IEEE, pages 73–78. <https://doi.org/10.1109/IMCSIT.2008.4747220>.
- Kalina Bontcheva, Hamish Cunningham, Ian Roberts, Angus Roberts, Valentin Tablan, Niraj Aswani, and Genevieve Gorrell. 2013. Gate teamwork: a web-based, collaborative text annotation framework. *Language Resources and Evaluation* 47(4):1007–1029.
- Bartosz Broda, Michał Marcińczuk, Marek Maziarz, Adam Radziszewski, and Adam Wardyński. 2012. KPWr: Towards a Free Corpus of Polish. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of LREC'12*. ELRA, Istanbul, Turkey.
- Hamish Cunningham, Diana Maynard, Kalina Bontcheva, Valentin Tablan, Niraj Aswani, Ian Roberts, Genevieve Gorrell, Adam Funk, Angus Roberts, Danica Damjanovic, Thomas Heitz, Mark A. Greenwood, Horacio Saggion, Johann Petrak, Yaoyong Li, and Wim Peters. 2011. *Text Processing with GATE (Version 6)*. <http://tinyurl.com/gatebook>.
- Richard Eckart de Castilho, Eva Mujdricza-Maydt, Seid Muhie Yimam, Silvana Hartmann, Iryna Gurevych, Anette Frank, and Chris Biemann. 2016. A web-based tool for the integrated annotation of semantic and syntactic structures. In *Proceedings of the workshop on Language Technology Resources and Tools for Digital Humanities (LT4DH) at COLING 2016*. pages 76–84.
- George Hripcsak and Adam S. Rothschild. 2005. Agreement, the f-measure, and reliability in information retrieval. *J. of Am. Medical Informatics Association* 12(3):296–298.
- Adam Kilgarriff, Vít Baisa, Jan Bušta, Miloš Jakubíček, Vojtěch Kovář, Jan Michelfeit, Pavel Rychlý, and Vít Suchomel. 2014. The sketch engine: ten years on. *Lexicography*.
- Jan Kocoń, Michał Marcińczuk, Marcin Oleksy, Tomasz Bernaś, and Michał Wolski. 2015. Temporal expressions in polish corpus kpwr. *Cognitive Studies—Études cognitives* (15):293–317.
- LCD. 2005. ACE (Automatic Content Extraction) English Annotation Guidelines for Events. Technical report, Linguistic Data Consortium.
- M. Marcińczuk and M. Ptak. 2012. *Preliminary study on automatic induction of rules for recognition of semantic relations between proper names in Polish texts*, volume 7499 LNAI.
- Michał Marcińczuk, Jan Kocoń, and Marcin Oleksy. 2017. *Liner2 — a generic framework for named entity recognition*. In *Proceedings of the 6th Workshop on Balto-Slavic Natural Language Processing*. Association for Computational Linguistics, Valencia, Spain, pages 86–91. <http://www.aclweb.org/anthology/W17-1413>.
- Michał Marcińczuk, Marcin Oleksy, Tomasz Bernaś, Jan Kocoń, and Michał Wolski. 2015. Towards an event annotated corpus of polish. *Cognitive Studies—Études cognitives* (15):253–267.
- Michał Marcińczuk, Michał Stanek, Maciej Piasecki, and Adam Musiał. 2011. Rich Set of Features for Proper Name Recognition in Polish Texts. In *SIIS 2011*. Springer.
- Michał Mirosław Marcińczuk, Marcin Oleksy, and Jan Wieczorek. 2016. Towards recognition of spatial relations between entities for polish. *Cognitive Studies—Études cognitives* (16):119–132.
- Małgorzata Marciniak, Agnieszka Mykowiecka, and Katarzyna Głowińska. 2010. Anotowany korpus dialogów telefonicznych. In Małgorzata Marciniak, editor, *Anotowany korpus dialogów telefonicznych*, Akademska Oficyna Wydawnicza EXIT, Warsaw, chapter Anotacja korpusu LUNA–WOZ.PL, pages 217–230.

- Michał Marcińczuk, Jan Kocoń, and Maciej Janicki. 2013. Liner2 – a customizable framework for proper names recognition for Polish. In Robert Bembenik, Lukasz Skonieczny, Henryk Rybinski, Marzena Kryszkiewicz, and Marek Niezgodka, editors, *Intelligent Tools for Building a Scientific Information Platform*, pages 231–253.
- Michał Marcińczuk, Monika Zaśko-Zielińska, and Maciej Piasecki. 2011. Structure annotation in the polish corpus of suicide notes. In Ivan Habernal and Václav Matoušek, editors, *Text, Speech and Dialogue*, Springer Berlin Heidelberg, volume 6836 of *Lecture Notes in Computer Science*, pages 419–426.
- Maciej Ogrodniczuk, Katarzyna Głowińska, Mateusz Kopeć, Agata Savary, and Magdalena Zawislawska. 2015. *Coreference in Polish: Annotation, Resolution and Evaluation*. Walter De Gruyter. <http://www.degruyter.com/view/product/428667>.
- Adam Radziszewski and Maciej Piasecki. 2010. A Preliminary Noun Phrase Chunker for Polish. *Proceedings of the Intelligent Information Systems* pages 169–180.
- Pavel Rychlý. 2007. Manatee/bonito - a modular corpus manager. In *1st Workshop on Recent Advances in Slavonic Natural Language Processing*. Masarykova univerzita, Brno, pages 65–70.
- M. Zaśko-Zielińska. 2013. *Listy pożegnalne: w poszukiwaniu lingwistycznych wyznaczników autentyczności tekstu*. Quaestio. <https://books.google.pl/books?id=QG60ngEACAAJ>.

INFOREX Corpora [Michał Marchinczuk] Logout

BSNLP 2017 shared task Agreement CCL Viewer Administration About & citing

Line2

Report 120364

Comparison

Only A	A and B	Only B
734940 [0,16]	Komisja Europejska	[bsnlp2017_org]
734322		
734844 [32,37]	Polskę	[bsnlp2017_loc]
734308		
734942 [72,91]	Puszczą Białowieżę	[bsnlp2017_loc]
734309		
734843 [92,108]	Komisja Europejska	[bsnlp2017_org]
734323		
734844 [137,156]	Puszczą Białowieżę	[bsnlp2017_loc]
734310		
734845 [221,228]	Polska	[bsnlp2017_loc]
734324		
734846 [292,301]	Natura 2000	[bsnlp2017_loc]
734325		
734847 [355,371]	Komisja Europejska	[bsnlp2017_org]
734326		
734848 [387,392]	Polskę	[bsnlp2017_loc]
734311		
734849 [427,446]	Puszczą Białowieżę	[bsnlp2017_loc]
734312		
734850 [457,461]	TOK FM	[bsnlp2017_org]
734313		
734851 [575,582]	Brukseli	[bsnlp2017_loc]
734314		
734852 [589,594]	Polska	[bsnlp2017_loc]
734315		
734853 [750,755]	Polskę	[bsnlp2017_loc]
734316		
734854 [758,794]	Europejskiego Trybunału Sprawiedliwości	[bsnlp2017_org]
734326		
734855 [938,956]	Planu Urządzenia Lasu	[bsnlp2017_misc]
734319		
734856 [958,969]	PUL	[bsnlp2017_misc]
734320		
734857 [965,986]	Nadlesnictwo Białowieża	[bsnlp2017_loc]
734328		
734858 [157,1578]	KE	[bsnlp2017_org]
734328		
734316 [292,301]	Natura 2000	[bsnlp2017_misc]
734317 [457,461]	TOK FM	[bsnlp2017_misc]
734318 [464,473]	TOKFM_NEWS	[bsnlp2017_misc]

View configuration

Annotation types

Annotation layer, subset or type

BSNLP 2017 shared task 4 (4) Display

Documents

By flag

names_min | gotowy

By subcorpus

IS Szydo Trump EC

Apply configuration

Agreement

Annotation category	Only A	A and B	Only B	PCS
bsnlp2017_org	12	319	20	95%
bsnlp2017_loc	10	132	5	95%
bsnlp2017_misc	12	59	29	74%
bsnlp2017_per	1	81	5	96%
all	35	591	59	93%

Involvement in development: Michał Marchinczuk, Adam Kaszmarek, Jan Kocot, Marcin Piak, Mikolaj Szewczyk, Marcin Oleksy, Grupa Technologii Językowych G4.19 Politechniki Wrocławskiej, 2009-2017

This page was tested in Firefox. Page generated in 0 sec(0).

Figure 5: Summary of annotation agreement for a set of document

Corpora
BSNLP 2017 shared task
Documents
Administration
About & citing

[Michał Marciniuk]
Logout

(41)
< poprzedni
> następny
+100
-10
ostatni >
(7)

Preview
Agreement
Metadata
Comment
Annotation
Annotation terms

Resolve annotations agreement

From	To	Text	User A	User B	Action for the final annotation
13	18	Polska	benlp2017_bc	benlp2017_bc	Add as benlp2017_bc (more)
20	37	European Commission	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
38	51	Unia Europejska	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
79	95	Funduszu Spójności	benlp2017_misc	benlp2017_misc	Add as benlp2017_misc (more)
254	255	M2	benlp2017_misc	.	Do not create an annotation (more)
					Add as benlp2017_misc (more)
409	421	Europejscykwo	benlp2017_per	benlp2017_per	Add as benlp2017_per (more)
455	486	Europejskiego Funduszu Społecznego	benlp2017_misc	benlp2017_misc	Add as benlp2017_misc (more)
488	504	Komisja Europejska	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
550	552	EFIS	benlp2017_misc	benlp2017_misc	Add as benlp2017_misc (more)
640	641	UE	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
643	648	Kanada	benlp2017_bc	benlp2017_bc	Add as benlp2017_bc (more)
678	690	Europejscykwo	benlp2017_per	benlp2017_per	Add as benlp2017_per (more)
721	724	CETA	benlp2017_misc	benlp2017_misc	Add as benlp2017_misc (more)
792	793	KE	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
795	800	Polsce	benlp2017_bc	benlp2017_bc	Add as benlp2017_bc (more)
802	815	Ambasada Kanady	benlp2017_org	benlp2017_org	Do not create an annotation (more)
810	815	Kanady	benlp2017_bc	benlp2017_bc	Do not create an annotation (more)
817	835	Ministerstwo Rozwoju	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
840	848	THINKTANK	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
870	873	CETA	benlp2017_misc	benlp2017_misc	Add as benlp2017_misc (more)
940	963	Mauro Raffaele Petitione	benlp2017_per	benlp2017_per	Add as benlp2017_per (more)
975	976	UE	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
980	983	CETA	benlp2017_misc	benlp2017_misc	Add as benlp2017_misc (more)
985	1001	Komisja Europejska	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
1003	1006	UNDP	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
1008	1028	Uniwerytet Warszawski	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
1033	1054	Fundacja Jozefa Robiata	benlp2017_org	benlp2017_org	Add as benlp2017_org (more)
1086	1106	Wykłady Kapsuśkiego	benlp2017_misc	benlp2017_misc	Add as benlp2017_misc (more)
1121	1160	Sali Balowej Pałacu Tyszkiewiczów-Potockich	benlp2017_bc	benlp2017_bc	Do not create an annotation (more)
					Add as benlp2017_bc (more)
					Add as benlp2017_bc (more)

Document content
View configuration

Strona główna | Polska - European Commission

Unia Europejska zamierza wydać ponad 422 mld euro z Funduszu Spójności w rozbudowę warszawskiego metra, z czego 100 mld przeznaczone jest na inwestycje w infrastrukturę kolejową. Dzięki inwestycji powstanie 6 nowych stacji na linii M2, terminal techniczny, do dotychczasowego labiru dołączy 13 nowych pociągów. Prace mają zostać zakończone do listopada 2019 r.

W okresie do 2007 do 2014 r. prawie dziesięciu milionom Europejczyków udało się znaleźć pracę dzięki pomocy z Europejskiego Funduszu Społecznego. Komisja Europejska opublikowała sprawozdanie oceniające inwestycje z EFS wraz ze szczegółowymi sprawozdaniami dotyczącymi poszczególnych krajów wspólnoty.

Unowa handlowa UE z Kanadą nadal wywołuje wątpliwości wśród Europejczyków? Czy polscy politycy sięczą z CETA? Jaką rolę w tym odegrał ministerstwo rozwoju? CETA w praktyce? CETA w praktyce? która Ambasada Kanady, Ministerstwo Rozwoju oraz THINKTANK zaprasza na sesję pt. "CETA w praktyce" która odbędzie się 28 stycznia 2017 r. Goszczem specjalnym będzie Mauro Raffaele Petitione, negocjator UE ds. CETA.

Komisja Europejska, UNDP, Uniwersytet Warszawski oraz Fundacja Jozefa Robiata zapraszają na obywaty wykład z cyklu Wykłady Kapsuśkiego. Już 17 stycznia w Sali Balowej Pałacu Tyszkiewiczów-Potockich przy ul. Krakowskie Przedmieście 32 w Warszawie będzie można wysłuchać wykładu wielokrotnie nominowanej do Pokojowej Nagrody Nobla Dawn Engle, twórczyni kampanii "One Billion Acts of Peace". Współorganizatorka i dyrektora fundacji PeaceJam – sieci zrzeszającej 13 zdobywców Pokojowej Nagrody Nobla - podzieli się swoim doświadczeniem w motywowaniu młodych ludzi do udziału na rzecz lokalnych społeczności.

Annotation layer, subset or type
Display

BSNLP 2017 shared task

4 (4)

Annotation name
Amns
A
B

Michał Marciniuk 35

Marcin Oleśny 34

Apply configuration

This page was tested in Firefox. Page generated in 0 sec(s).

Involved in development: Michał Marciniuk, Adam Kaczmarek, Jan Kocoh, Marcin Paik, Mikołaj Szewczyk, Marcin Oleśny, Grupa Technologii Językowych G4-19 Politechniki Wrocławskiej, 2009–2017

Figure 6: User agreement verification for a single document