

Semantic annotation and expansion for keyword queries part2, hands-on

Ranka Stanković
ranka@rgf.rs

<http://rgf.rs/ranka2.pptx>

University of Belgrade, Serbia



COST Action IC1302

2nd KEYSTONE Training School

Keyword search in Big Linked Data



Ranka Stanković. "Semantic annotation and expansion for keyword queries". The 2nd KEYSTONE Training School on Keyword search in Big Linked Data, University of Santiago de Compostela, Spain, 18-22 July 2016



Main topics

- Tools and guidelines
- Tool for enhanced search of bilingual digital libraries
- Examples of application in GIS
- Linguistic analysis and annotation
 - Math
 - Unitex
 - WebCQP
- Where To Go From Here?

Ontotex - GraphDB

Free RDF triplestore

- Semantically Rich NoSQL Graph Database
- data management for very large sets of structured, semi-structured or unstructured data.
- NoSQL ('not only SQL') graph databases serve organizations to access, integrate and analyze both unstructured data and data stored in the cloud, thus helping them with their big data and social media analytics.
- More on <http://ontotext.com/knowledge-hub/>
- Any experiance?
- Discussion?
- Beter option?

W3C Multilingual Linked Open Data Community Group

- Best Practices for Multilingual Linked Open Data Community Group (December 2015)

<https://www.w3.org/community/bpmlod/>

- Guidelines for Linguistic Linked Data Generation: Multilingual Dictionaries (BabelNet)
- Guidelines for Linguistic Linked Data Generation: Bilingual Dictionaries
- Guidelines for Linguistic Linked Data Generation: Multilingual Terminologies (TBX, TermBase eXchange)
- Guidelines for developing NIF-based NLP Web Services
- Guidelines for LLD Exploitation

NLP Interchange Format (NIF)

- NIF is an RDF/OWL-based format that aims to achieve interoperability between NLP tools language resources and annotations
- best practices to follow for the implementation of RESTful NLP web services that rely on NIF
- They implemented NIF wrappers for the *Stanford POS tagger* and *Stanford parser*
- <http://www.w3.org/2015/09/bpmlod-reports/nif-based-nlp-webservices/>

Parallel Corpora

- United Nations Parallel Corpus

<http://conferences.unite.un.org/UNCorpus>

 **OPUS** open parallel corpus search & browse

- <http://opus.lingfil.uu.se/>
 - OPUS multilingual search interface
 - Europarl v7 search interface
 - Europarl v3 search interface
 - OpenSubtitles search interface
 - EUconst search interface
 - Word Alignment Database

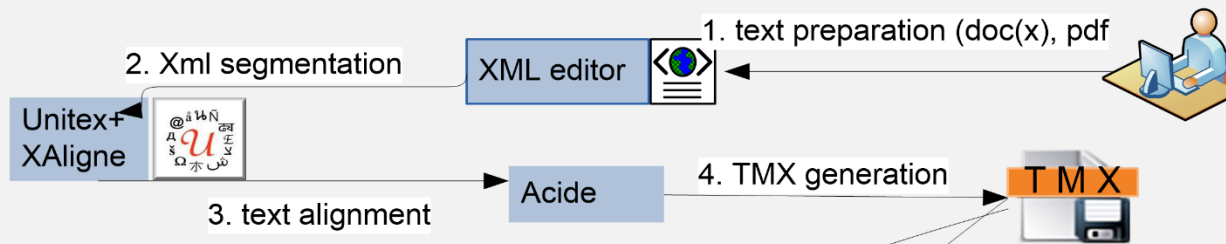
Task 1

- Find and select parallel corpus on offered links
- Search & browse
- Find another (not listed on slides) in your country / region / language
- Compare: size, user interface, content

Query expansion for aligned text

- Bibliša (<http://jerteh.rs/Biblisha>) is a tool for enhancement of search possibilities in multilingual digital libraries of e-journals
 - Cross-lingual search functions for large collections of aligned texts
 - Queries with simple and multiword keywords in Serbian and English
 - Queries expanded semantically and morphologically
- Previously enhanced cross-lingual search implemented in LeXimir with XPath queries works only for one document
- Extension of cross-lingual search for collections relies on XML databases and XQuery language

Text alignment and TMX document creation

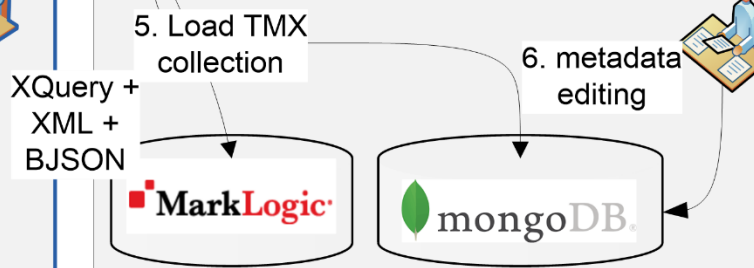


Keyword based search

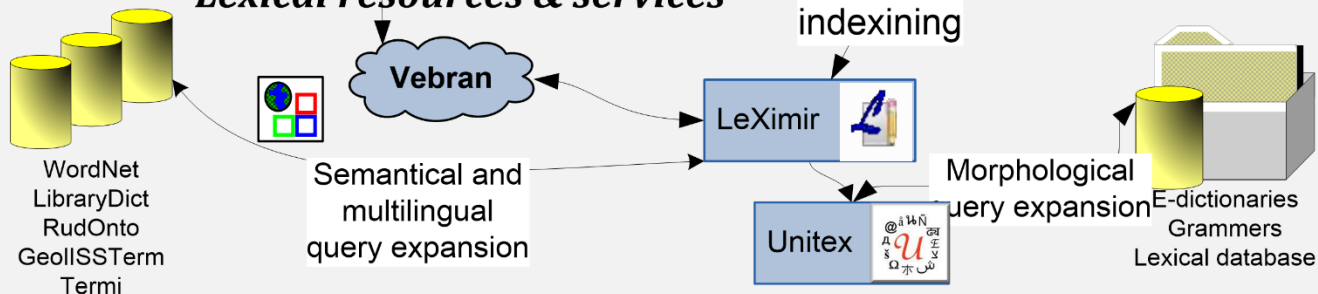


Keyword query expansion

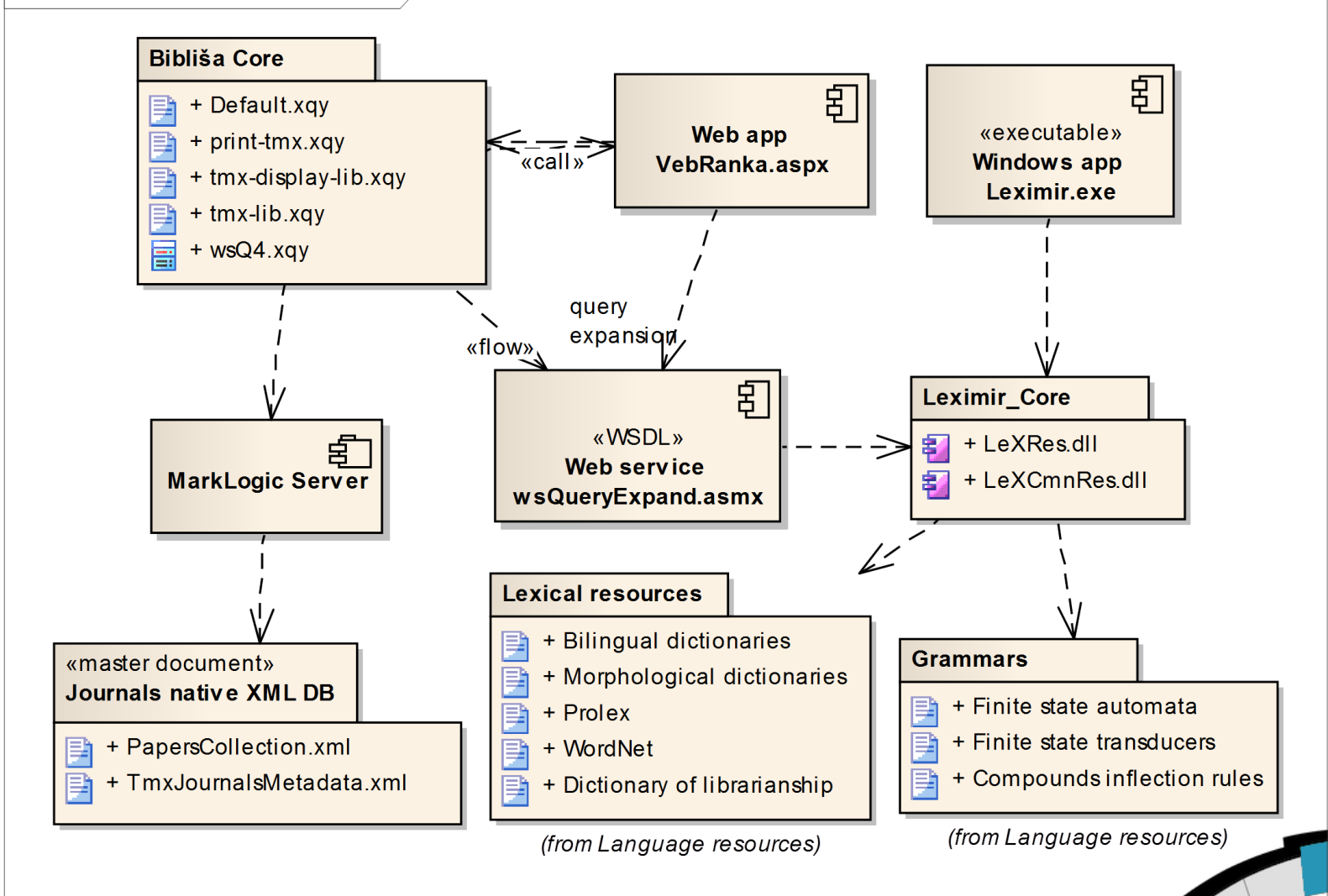
Bilingual digital library



Lexical resources & services



cmp Bibliša+Leximir_Components



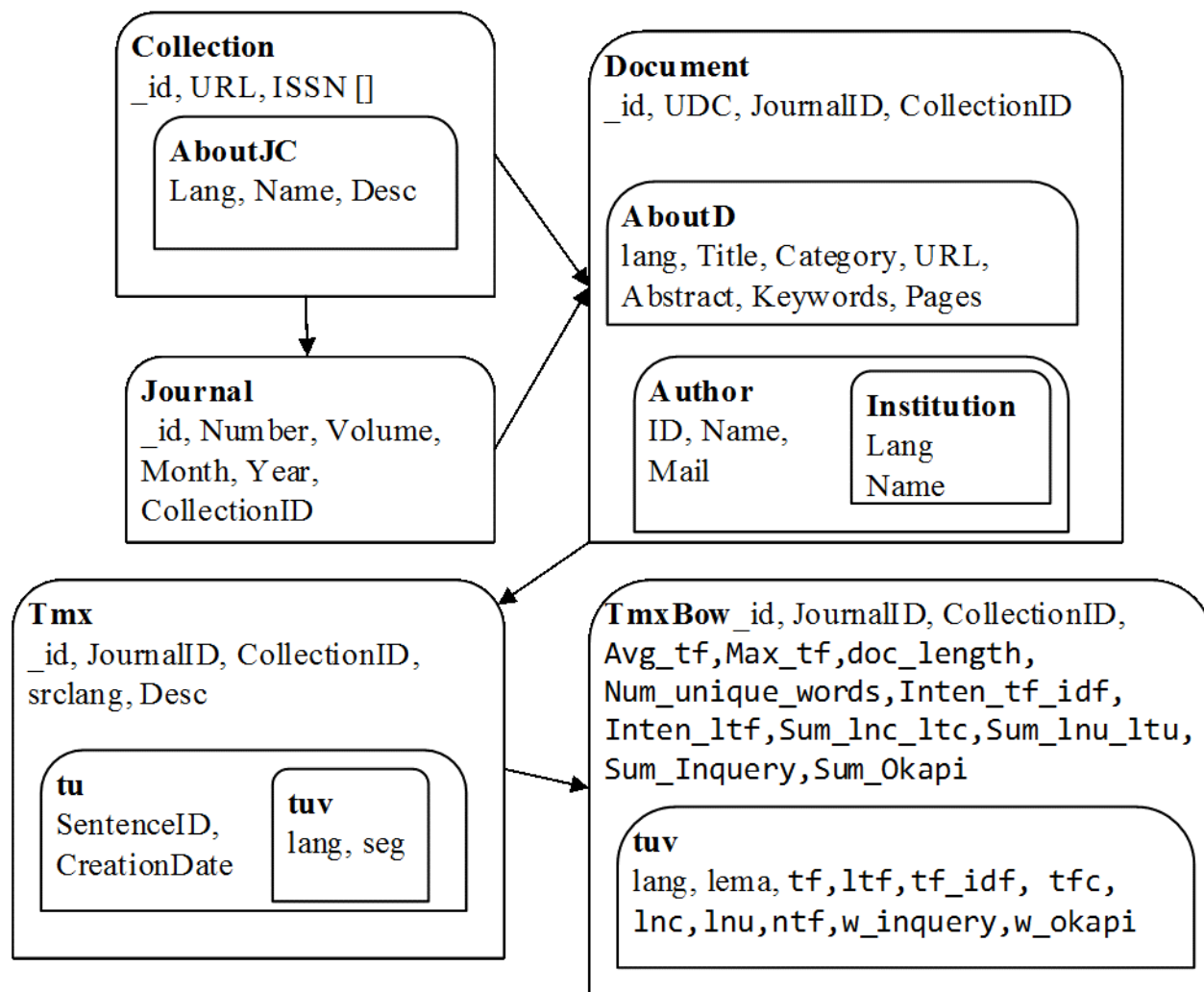
Resources

```
<tu>
  <prop type="Domain">Tomašević et al., 2012, vol. XX:20, ID: 2.2012.20.4
  </prop>
  <tuv xml:lang="en" creationid="n5 " creationdate="20140525T090842Z">
    <seg>It is the authors' wish to bring closer to the mining public at
    least a small part of the possibilities offered by GIS, both in mining
    activities and in environment protection. </seg>
  </tuv>
  <tuv xml:lang="sr" creationid="n5 " creationdate="20140525T090842Z">
    <seg>Želja autora je da rudarskoj javnosti približi samo mali deo
    ogromnih mogućnosti koje GIS pruža, kako pri rudarskim aktivnostima,
    tako i u zaštiti životne sredine. </seg>
  </tuv>
</tu>
{"Desc": "Tomašević et al., 2012, vol. XX:20, ID: 2.2012.20.4",
 "SentenceID": "n5",
 "tuv": [{
   "lang": "en",
   "seg": "It is the authors' wish to bring closer to the mining
   public at least a small part of the possibilities offered by
   GIS, both in mining activities and in environment protection. "
 },
 {
   "lang": "none",
   "seg": "Želja autora je da rudarskoj javnosti približi samo mali
   deo ogromnih mogućnosti koje GIS pruža, kako pri rudarskim
   aktivnostima, tako i u zaštiti životne sredine. "
 }
 ]}
```

- MarkLogic (<http://www.marklogic.com>) server for management of document collections consisting of aligned parallel texts converted in TMX (Translation Memory eXchange) format and MongoDB (recently)
- Collection of TMX documents generated by ACIDE with export routines for automatic generation of TMX versions of parallel texts aligned by Xalign



Resources - metadata



Software implementation

- Excerpt from XQuery code: variable \$x creates a collection of “tu” XML nodes (elements), which can subsequently be queried, yielding the results obtained by \$query tagged with :

```
for $x in cts:search(fn:doc()//tu,$query)
return
cts:highlight($x, $query,<em>{$cts:text}</em>)
```

- The results, a set of aligned concordances, are formatted and presented to the user
- The concordances are preceded by information identifying the document they originate from, and a link to summary metadata for this document in both languages

Supporting resources

- Four types of resources are used for the expansion of queries submitted to our collection of documents.
- Serbian morphological dictionaries of simple words and multi-word units developed within: Unitex and Nooj. are used to generate all inflective forms of query keywords, thus improving the system recall without negative effects on precision: 135 000 lemmas of simple words and 20 000 compounds.
- Grammars in the form of finite state automata and transducers, amended by compounds inflection rules
- Princeton English Wordnet (PWN), version 2.0 and Serbian Wordnet (SrpWN), initially developed in the scope of the BalkaNet project and subsequently enhanced and upgraded, 17 500 synsets with 30 000 literals
- Dictionary of librarianship: English-Serbian and Serbian-English 23 400 terms – 11 300 in English and 12 100 in Serbia

```

<request>
  <query xml:lang='sr'>digitalna biblioteka</query>
  <query xml:lang='sr'>digitalnih biblioteka</query>
  <query xml:lang='sr'>digitalne biblioteke</query>
  <query xml:lang='sr'>digitalnim bibliotekama</query>
  ....
  <query xml:lang='en'>digital library</query>
</request>

```

```

cts:or-query ((cts:word-query ("digitalna biblioteka", ("stemmed", "lang= sr")) ).
cts:or-query ((cts:word-query ("digitalnih biblioteka", ("stemmed", "lang= sr")) ).
cts:or-query ((cts:word-query ("digitalne biblioteke", ("stemmed", "lang= sr")) ),
cts:or-query ((cts:word-query ("digitalnim bibliotekama", ("stemmed", "lang= sr")) ),
.....
cts:word-query (digital library, ("stemmed", "lang= sr") ) )

```


http://hlt.rgf.bg.ac.rs:8005/wsQ4.xqy?request=<request><query%20xml:lang= Live Search		
Biblisha concordances		
Đukić, 2011, vol. XII:1, ID: 2011.1.10 metadata	n21 : The AccessIT project seeks to deliver a unique package of practical training and skills development, supported by clear guidance, to enable smaller, local cultural organizations in countries where progress in this area is currently limited, to maximize the opportunities provided by the new technologies (combined with major policy implementations such as the <i>Digital libraries Initiative</i>) to deliver and disseminate arts and cultural offerings to the citizens of Europe most effectively.	n21 : Projekat AccessIT je obezbedio jedinstveni paket obuke i usavršavanja za osposobljavanje stručnjaka iz manjih, lokalnih ustanova kulture - baštinskih institucija iz Turske, Grčke i Srbije za digitalizaciju fondova biblioteka, arhiva i muzeja.
Đukić, 2011, vol. XII:1, ID: 2011.1.10 metadata	n29 : One of the expected results of the project is creation of national repositories (or at least making necessary conditions for it) which will be able to join European digitization projects (such as Europeana and Europeana Local).	n29 : Kao jedan od rezultata projekta očekuje se stvaranje nacionalnih repozitorijuma (ili barem uslova za njihovo stvaranje) koji će biti spremni da se uključe u projekat Evropske <i>digitalne biblioteke</i> (Europeana i Europeana Local).
Đukić, 2011, vol. XII:1, ID: 2011.1.10 metadata	n32 : On the basis of the cooperation in this project BCI signed agreement with PSNC which allows BCI to use dLibra software to create its own <i>digital library</i> .	n32 : Na osnovu saradnje u okviru ovog projekta Biblioteka grada Beograda je sa PSNC potpisala poseban sporazum o saradnji na osnovu koga je dobila mogućnost korišćenja softvera dLibra za potrebe izgradnje sopstvene <i>digitalne biblioteke</i> .
Filipi Matutinović, 2011, vol. XII:1, ID: 2011.1.8 metadata	n4 : The most well-known is Google Books, which started as a private company project with the idea to build world <i>digital library</i> by scanning books from biggest public and university libraries.	n4 : Najpoznatiji je Google Books, koji je započet kao projekat privatne kompanije sa idejom da se izgradi svetska <i>digitalna biblioteka</i> putem skeniranja knjiga iz najvećih javnih i univerzitetskih biblioteka.
Filipi Matutinović, 2011, vol. XII:1, ID: 2011.1.8 metadata	n8 : European Commission proclaimed Digital Agenda and in 2008 launched the project Europeana, European <i>digital library</i> with 2 million objects from EU countries' museums, archives and libraries.	n8 : Evropska Komisija je proklamovala Digitalnu agendu i započela projekat Europeana, Evropsku <i>digitalnu biblioteku</i> sa 2 miliona objekata iz muzeja, arhiva i biblioteka zemalja Evropske unije.
Butigan-Vučaj, 2011, vol. XII:1, ID: 2011.1.9 metadata	n3 : The meaning of the acronym of the EMBARK project is associated with the main goal of the project, which is "embarking" of certain number of Cyrillic manuscript books into the Manuscriptorium - European <i>digital library</i> of written historical resources.	n3 : Akronim projekta EMBARK na engleskom jeziku znači ukrcavanje. A cilj projekta jeste "ukrcavanje" jednog broja ćirilskih rukopisa u Manuscriptorium - evropsku <i>digitalnu biblioteku</i> pisanih istorijskih izvora.

Main features of *Bibliša*

- A complex system composed of several modules
- Targeted at textual resources in the form of collections of TMX documents and the corresponding metadata
- Able to expand search queries:
 - morphologically, using morphological e-dictionaries, the system of rules for compound inflection, and finite automata and transducers
 - semantically and bilingually, using Serbian and English wordnets and the bilingual Dictionary of Librarianship

E-dictionaries	Serbian wordnet	Dictionary of Librarianship
128,000 simple word lemmas	17,500 synsets	11,300 English terms
10,000 compound word lemmas	30,000 literals	12,100 Serbian terms



Query realization in *Bibliša*

- The user formulates the initial query as one or more keywords (simple or multiword)
- If the user so specifies, *Bibliša* forwards this query for further morphological and semantic expansion
- The expansion is handled by a web service that is part of the LeXimir software package, a multipurpose tool also developed by the HLT Group
- The web service invokes LeXimir's function library LeXimirCore, whose functions expand the query, using available lexical resources and Unitex routines
- The expanded query is transformed into an XQuery and used for searching the TMX document collection obtained from journal articles
- As a result a set of aligned concordances is obtained and presented to the user



What we can learn?

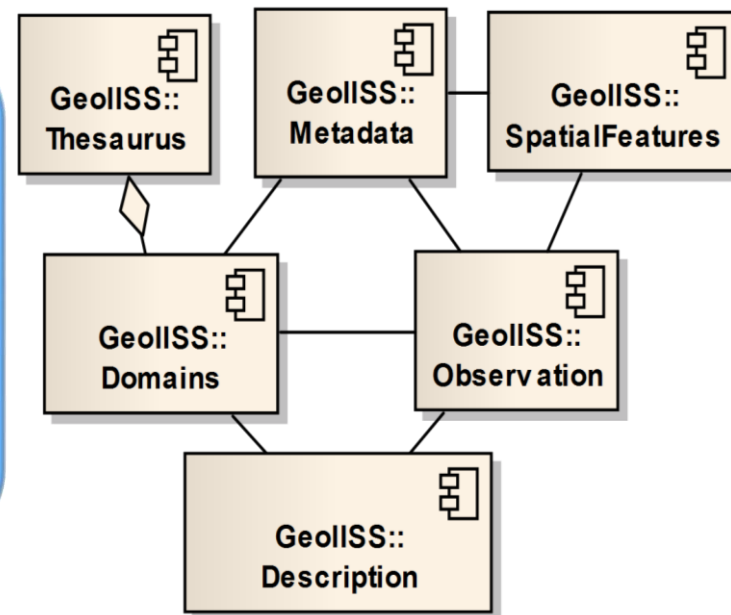
- Check context of terms
- Check other options for translation (partial match)
- Expand dictionary (with partial match)
- Use TMX for word alignment
- Aligned word for query expansion
- Machine translation

Task 2 Query expand using Biblisha

- <http://jerteh.rs/biblisha/>
- Metadata search
- Fulltext search
- Expand query (have in mind domains of resources)
- Search in Mongo and Mark Logic
- Differences?
- MongoDB multilingual search experience?
- Best DBMS solution for multilingual parallel corpora?

Query expansion in GIS

Stanković Ranka, Obradović Ivan, Kitanović Olivera, "GIS Application Improvement with Multilingual Lexical and Terminological Resources", Proceedings of the 5th International Conference on Language Resources and Evaluation, LREC 2010, Valetta, Malta, May 2010 , pp.2283-2287. European Language Resources Association Valetta, Malta 2010 ISBN:2-9517408-6-7 [\[link\]](#)

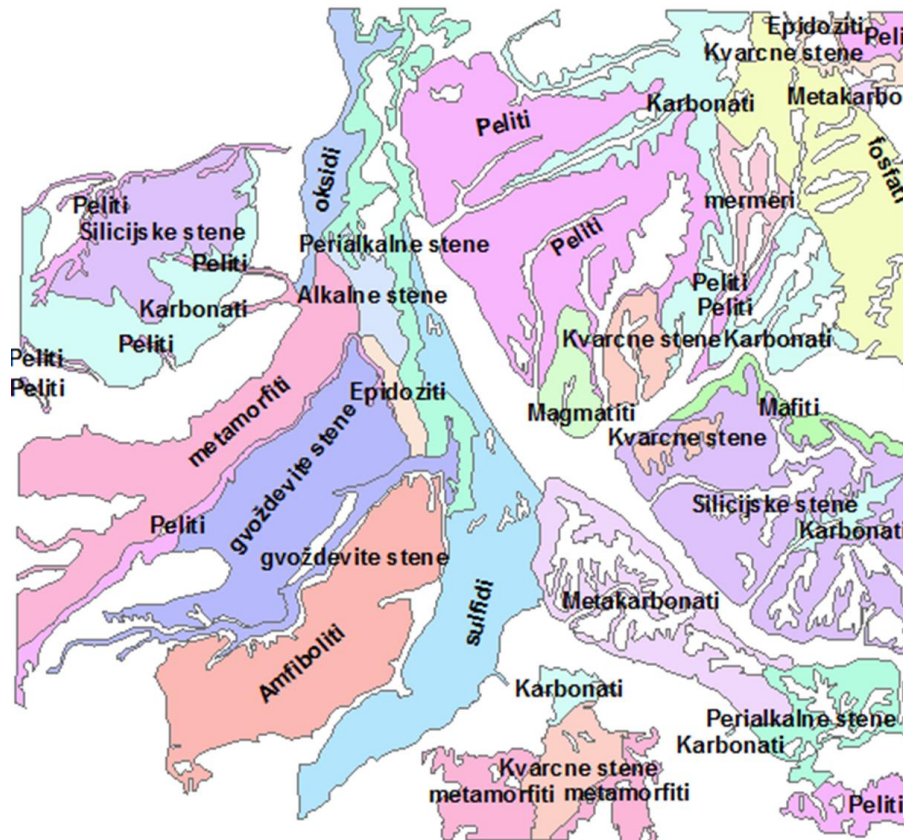


- GeolISS

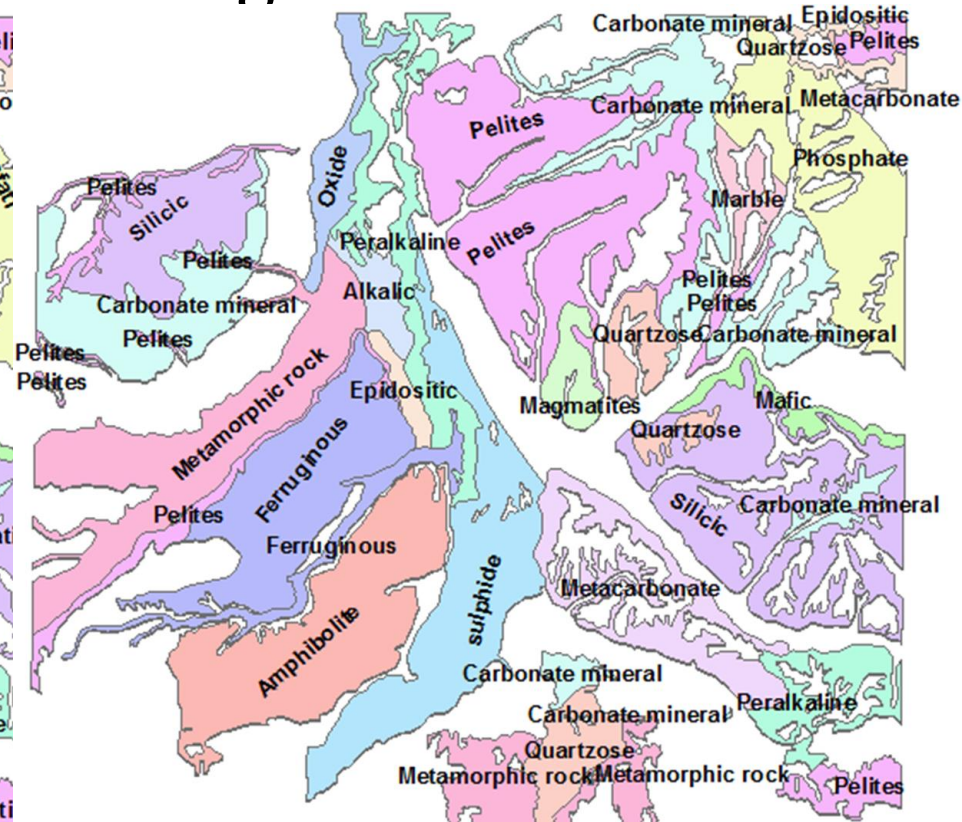
- Digital archiving, query, retrieving, analysis and visualization of geological data
- Integration of existing geologic archives, data from published maps, newly acquired field data
- Domains controlled by the geologic thesaurus
- Implemented using .NET + ESRI ArcGIS

Multilingual search and annotation

Serbian



English



```

Synset: zemljotres:1, potres:1b, trus:1
RevTree  Unitex Graph  Text  HH Tree  XML
├─ [n] fenomen:1, pojava:1a
│   └─ H*[n] prirodna pojava:X, prirodni fenomen:X
│       └─ H[n] geoloska pojava:1
│           └─ H[n] zemljotres:1, potres:1b, trus:1
│               └─ H[n] udar:5, talas:1
│                   └─ H[n] podhtavanje tla:2, slab zemljotres:1
│                       └─ H[n] podvodni zemljotres:1
  
```

QE by WordNet

Serbian & English WordNet
HH tree for synsets from
domain 'geology

- Other HLT Group resources
 - morphological e-dictionaries in LADL format
 - finite state transducers developed within the Unitex system
 - electronic thesauri
 - ontologies
- HLT Group applications
 - LeXimir (various language processing related tasks)
 - Web Service for Query Expansion

```

Synset: earthquake:1, quake:1, temblor:1...
SEM mark  RevTree  Unitex Graph  Text  HH Tree  XML
phenomenon:1
H*[n] natural phenomenon:1
├─ H[category_domain][n] geological phenomenon:1
│   └─ H[eng_derivative][n] earthquake:1, quake:1, temblor:1, seism:1
│       └─ H[eng_derivative][eng_derivative][v] tremor:1, quake:2
│           └─ H[n] shock:5, seismic disturbance:1
│               └─ H[eng_derivative][n] tremor:2, earth tremor:1, microseism:1
│                   └─ H[n] seaquake:1, submarine earthquake:1
  
```



voz sadrži počinje sa tačna fraza
 Literal Def Usage Domain

[Sinsetovi korisnika](#)

Ukupno nađeno: **1** sinset

ID: ENG30-04468005-n POS: n BCS: 3 0.000 0.000 User
 21.07.2004 Approved: **yes** PWN XML

Literals: **voz (1), vlak (1)**
 Definition: *Javni prevoz koji obezbeđuje kompozicija vagona koje vuče lokomotiva.*

- ▼ - Relations... [hyperym-> ENG30-04019101-n, sredstvo javnog prevoza](#)
- ▼ - Relations... [eng_derivative-> ENG30-01936537-v, putovati vozom](#)
- ▼ - Relations... [mero_member-> ENG30-02959942-n, vagon, kola](#)
- ▼ - Relations... [mero_member-> ENG30-03684823-n, mašina, lokomotiva](#)

SUMO: Train =
 DOMAIN: transport

train Word Sinset ID
 Tree View

Number of Nouns: 6

ID {3431745} Sense {{gearing, gear, geartrain, power_train, train}: wheelwork consisting of a connected set of rotating gears by which force is transmitted or motion or torque is changed; "the fool got his tie caught in the geartrain"}

sumo: { Device + } domain: {mechanics} pos: {0}, neg: {0} SWN

▼ - Relations...

ID {4468476} Sense {{train}: piece of cloth forming the long back section of a gown that is drawn along the floor; "the bride's train was carried by her two young nephews"}

sumo: { Clothing + } domain: {fashion} pos: {0}, neg: {0} SWN

▼ - Relations...

ID {7294777} Sense {{train}: a series of consequences wrought by an event; "it led to a train of disasters"}

sumo: { result + } domain: {} pos: {0}, neg: {0} SWN

▼ - Relations...

ID {8427629} Sense {{caravan, train, wagon_train}: a procession (of wagons or mules or camels) traveling together in single file; "we were part of a caravan of almost a thousand camels"; "they joined the wagon train for safety"}

sumo: { Transportation + } domain: {transport} pos: {0}, neg: {0} SWN

▼ - Relations...

ID {8459648} Sense {{string, train}: a sequentially ordered set of things or events or ideas in which each successive member is related to the preceding; "a string of islands"; "train of mourners"; "a train of thought"}

sumo: { Collection + } domain: {} pos: {0}, neg: {0} SWN

▼ - Relations...

ID {4468005} Sense {{train, railroad_train}: public transport provided by a line of railway cars coupled together and drawn by a locomotive; "express trains don't stop at Princeton Junction"}

sumo: { Train = } domain: {transport} pos: {0}, neg: {0} SWN

▼ - Relations...

Serbian semantical resources

<http://sm.jerteh.rs>

Task 3. Search Wordnet

- Search Serbian WN (serbian left side, english right side)
- Search Princeton on their site
- Find WN in your (favorite) language
- Compare?
- Suggestions?
- Best solution?

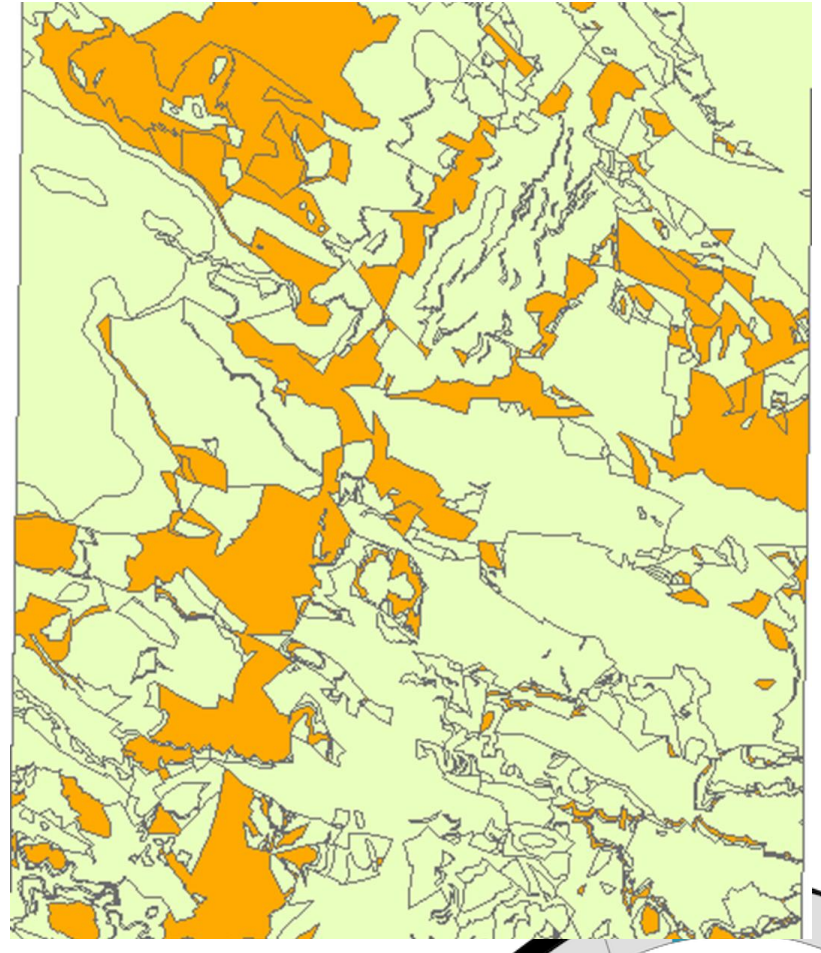
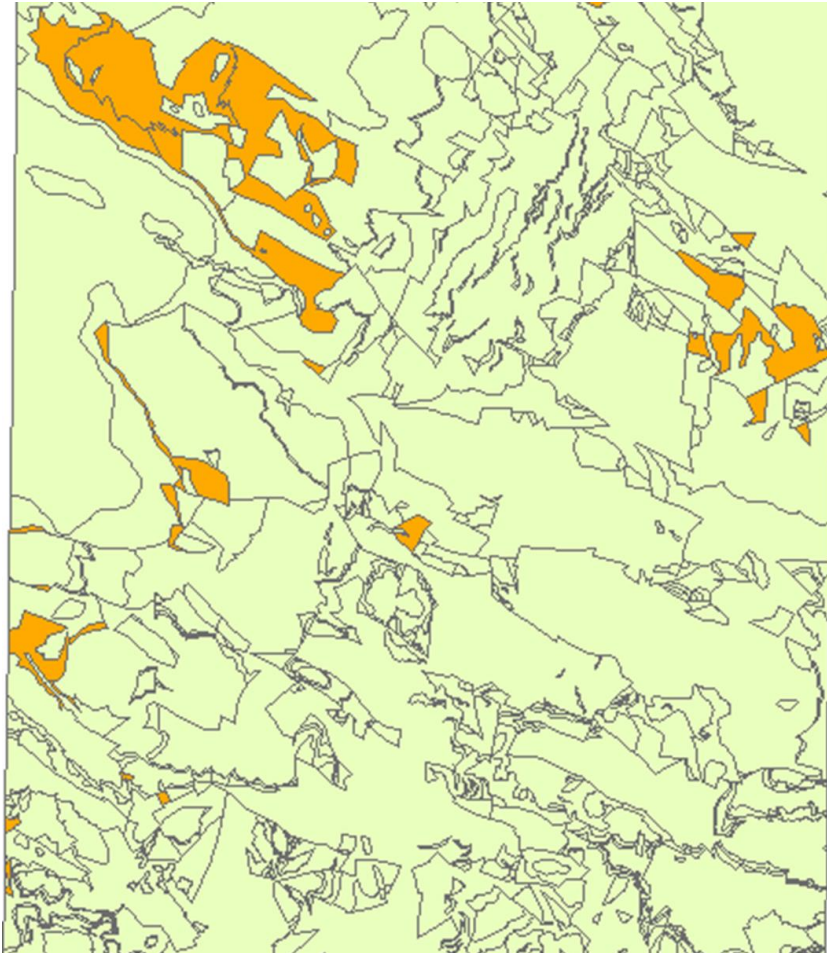
QE by Geological thesaurus

- GeolISSTerm, aggregation of geological vocabularies
 - {ID=2356, Name=Piroklastična stena, Def= Stena nastala depozicijom i litifikacijom piroklastičnih naslaga; fragmenti ove stene obrazovani su direktnom eksplozivnom fragmentacijom, Synonym= Piroklastit, Hyperonym= 2123}
- multilingual
 - {ID=12345, Name=Pyroclastic rock, Def=A volcanoclastic rock formed by direct explosive volcanic activity, OrigID=2356, Lng=EN)
 - {ID=12367, Name= Roches pyroclastiques, Def= Une pierre volcanoclastiques formées par l'activité volcanique explosive direct, OrigID=2356, Lng=FR)

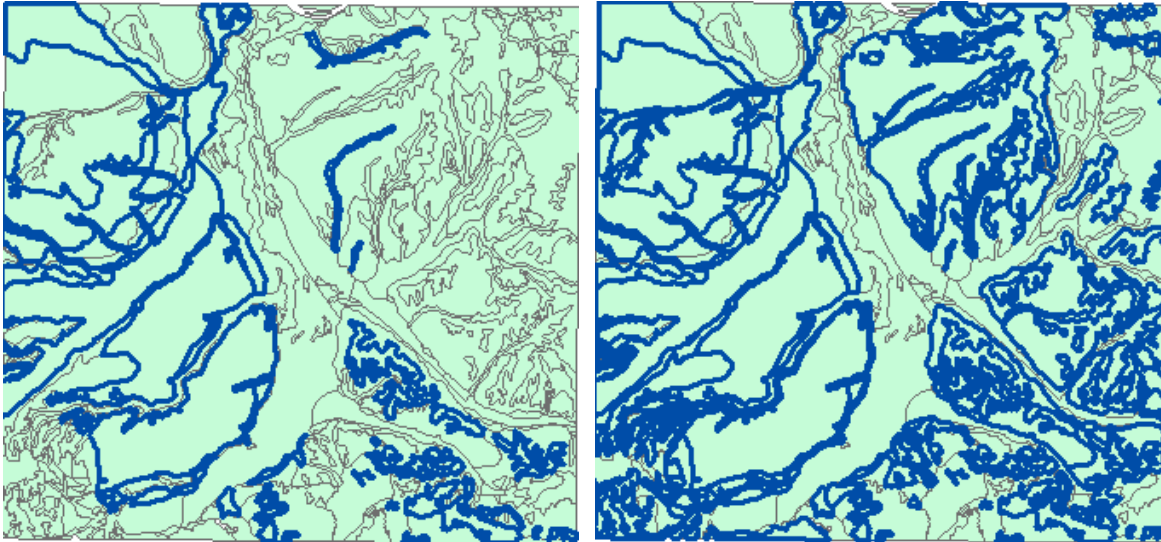
GeolISS Query Expansion

- Selection of keywords can be improved by using lexical resources, morphological dictionaries and transducers
- Morphological expansion
 - retrieval of geological units containing 'limestone' in their description field ('krečnjak' in Serbian)
>>> 95
 - with morph.exp.: 'krečnjak, krečnjaka, krečnjaku, krečnjakom, krečnjače, krečnjaci, krečnjacima, krečnjake'
>>> 249
- GeolISSTerm Semantic expansion
 - 'Middle Jurassic' (srednja jura) >>> 53 + doger (dogger) >>> 134
 - Lithologic constituent (+ Lithologic Member + Lithologic 'Package') no improvement
 - WordNet Semantic expansion 'izvor, vrelo' (eng. 'spring, outflow, outpouring, natural spring')

Selected map units with the original and morphologically expanded query



QE in GIS



Morphological expansion

glina (clay) 51

→

glina, glinama, glinom, glinu, glini, gline 118



Semantic expansion

Lower Pliocene (donji pliocen) 13

→

donji pliocen, pont 33



What is EuDML?

EuDML makes the mathematics literature available online in the form of an enduring digital collection, developed and maintained by a network of institutions.

<https://eudml.org/search/>

The screenshot shows the EuDML website's search interface. At the top, there is a navigation bar with tabs for Home, Advanced Search, Browse by Subject, Browse by Journals, and Refs Lookup. The Advanced Search tab is active. Below the navigation bar, there is a search bar with the text 'Title, Author, Keyword, Citation, Date...' and a search button. The search results are displayed in a table with columns for 'Currently displaying' and 'Languages'. The search criteria are 'Any field contains circle' and 'Contains the following math formula (red border means the formula is incomplete)'. The formula is $x^2 + y^2 = 1$. The search results show a list of documents, including 'The polynomial $x^3 + x^2 + x - 1$ and elliptic curves of conductor 11' by Alfred J. Van der Poorten.

English (en) | Login | Register | (Why Register?)
Title, Author, Keyword, Citation, Date... Search

Home | **Advanced Search** | Browse by Subject | Browse by Journals | Refs Lookup

Advanced Search [Back to Simple Search](#)

Match of the following rules

Any field contains [Add Sub-clause](#) [Add Another Rule](#)

Contains the following math formula (red border means the formula is incomplete)

Formula preview
 $x^2 + y^2 = 1$

Only documents with accessible full-text

Search

Currently displaying 1 – 20 of 4802
Showing per page
Order by [Relevance](#) | [Title](#) | [Year of publication](#)

[The polynomial \$x^3 + x^2 + x - 1\$ and elliptic curves of conductor 11](#)
Alfred J. Van der Poorten
Séminaire Delange-Pisot-Poitou. Théorie des nombres

Languages

en	4390
fr	187
cs	77
es	16
de	15



Task 4. Search EuDML

- Find documents for circle with or without formula
- Change names of variables. Compare results.
- Filter language
- Same on other library (math)
- Find new (not listed)
- Best solution?

Unitex

Demo for english, french, serbian,...

- <http://www-igm.univ-mlv.fr/~unitex>

Home
[Why Unitex ?](#)
[Screenshots](#)
[Download](#)
[User manual](#)
[Forum](#)
[Bug Reporting Guide](#)
[Language resources](#)
[LGPLLR licensed data](#)
[LGPL](#)
[LGPLLR](#)
[Your contribution](#)
[Links](#)
[Bibliography](#)
[Works with Unitex](#)
[Mailing list](#)
[Unitex Library - User's Guide](#)
[Student Project Proposals](#)

Unitex/GramLab is an open source, cross-platform, multilingual, lexicon- and grammar-based corpus processing suite

[Unitex/GramLab 3.1 Stable is now available](#)

Unitex/GramLab has been selected as a Google Summer of Code 2016 mentor organization

Google Summer of Code (GSoC) is a global program that offers students stipends to write code for open source projects during summer break. This year, Unitex/GramLab has been selected as a Google Summer of Code mentor organization. If you're interested in helping with GSoC, mentoring a student, or you are a student, we'd love to hear from you:

- [Our organisation profile](#)
- [View our ideas list](#)
- [Google Summer of Code 2016 website](#)
- [More organisations in the "languages" category](#)
- [GSoC how it works](#)
- [GSoC timeline](#)
- [GSoC FAQ](#)

If you have any questions, please do not hesitate to post back at the users [forum](#) or to send a message to the [developers mailing list](#)

On the Unitex/GramLab forum, you can ask and answer questions and post your suggestions about Unitex and GramLab.

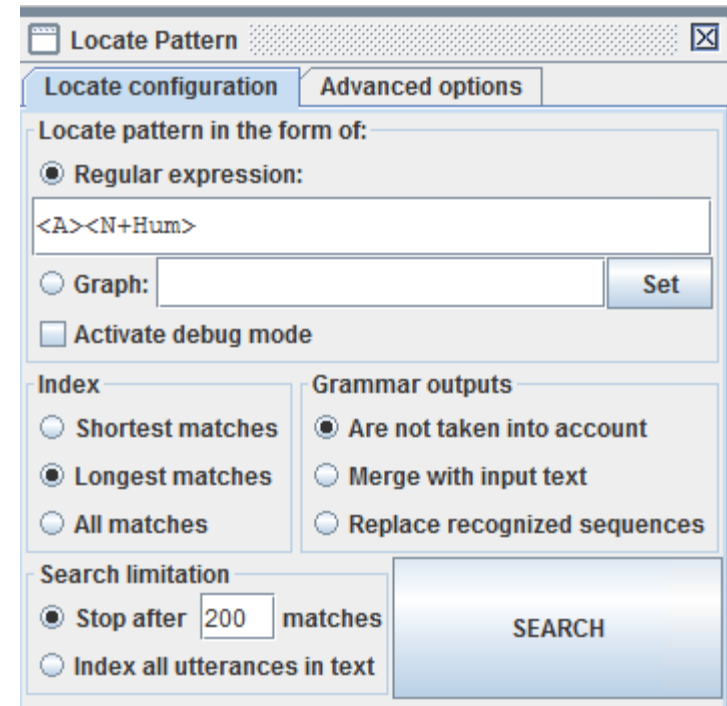
[Unitex/GramLab Forum](#)

According to a study based on 377 job offers for NLP engineers from March 2013 to July 2015, Unitex is among the **most expected skills in terms of NLP tools (9%)**

Task 5a. Unites

- Run
- Select language
- Open text document (Text->Open)
- Process (automatic, default)
- Text -> Locate pattern (fig.)
- Build concordances

Paumier, S. "Unitex 3.1 User Manual (2016).", <http://www-igm.univ-mlv.fr/~unitex/>



Task 5b. Query in UniteX

- **<to be>** - all word forms linked to lemma
- **<N>** - search for POS – all nouns
- **<DET><A><N+Hum>** - noun preceded by determiner and adjective
- **<N+Hum>** - nouns with semantic tag Human
- **<N+NProp+Hum~Inh>** - ...not Inhabitans
- **(<A>+<PRO+ProA>) <love>**
- **(<this>+<that>) <A> <N>**
- **Try yours idea!**

Task 6 CQPweb

- Try
 - <http://cwb.sourceforge.net/cqpweb.php>
 - <https://cqpweb.lancs.ac.uk>
- CQP language
 - http://cwb.sourceforge.net/files/CQP_Tutorial/
- Find other
- Favorite?

Thank you for your
attention

Hvala na pažnji
Хвала на пажњи