

E-Connecting Balkan Languages

Cvetana Krstev
Faculty of Philology
University of Belgrade
cvetana@matf.bg.ac.rs

Ranka Stanković
Faculty of Mining and
Geology
University of Belgrade
ranka@rgf.bg.ac.rs

Duško Vitas
Faculty of Mathematics
University of Belgrade
vitas@matf.bg.ac.rs

Svetla Koeva
Dep. of Computational
Linguistics
Institute for Bulgarian
svetla@dcl.bas.bg

Abstract

In this paper we present a versatile language processing tool that can be successfully used for many Balkan languages. For its work, this tool relies on several sophisticated textual and lexical resources that have been developed for most Balkan languages. These resources are based on several *de facto* standards in natural language processing.

Keywords

Query expansion, e-dictionaries, wordnets, proper names, aligned texts

1. Introduction

The software tool WS4LR (shortened for WorkStation for Language Resources) has been developed by the Language Technology Group organized at the Faculty of Mathematics for several years now. Its first version was introduced in 2004 [8] and it dealt mainly with harmonizing various heterogeneous lexical resources. Subsequently, many new features have been added, particularly those that have helped in the production and exploration of aligned texts on the basis of the lexical resources incorporated [9]. The new tool WS4QE (shortened for Work Station for Query Expansion) was developed on the basis of WS4LR that enables the expansion of queries submitted to the Google search engine [10]. The integrated lexical resources enable modifications of users' queries for both monolingual and multi lingual searches.

When presenting WS4LR and WS4QE, we have always stressed that, although they have been mainly used for Serbian, they are by no means language dependent as long as compatible lexical resources exist for any two languages. Nevertheless, the full potential of these tools was until now used only for Serbian, and in bilingual context, for Serbian and English.

In this paper we will show that the tools WS4LR and WS4QE are truly independent both from Serbian, for which they were initially developed, and from English which seems to be in the background of many natural language processing tools. The main presupposition for the usage of these tools for other languages is the existence of textual and lexical resources developed in the same methodological framework. Since this prerequisite has been satisfied for Bulgarian, and, to some extent, for some other Balkan languages (Greek, Romanian, etc...), we will

therefore show that WS4LR and WS4QE can be successfully used for these languages.

WS4LR, written in C#, is organized in modules which perform different functions. The core of the system WS4LR_Core comprises four .Net libraries, used by two components: the stand-alone windows application WS4LR.exe and the web service wsQueryExpand.asm. Web application WS4QE.asp manages user query request, than uses web service in order to expand user query, submits the expanded query to Google search engine and finally presents retrieved result.

2. Integrated Language Resources

In order to prove the usability of WS4LR and WS4QE for languages other than Serbian and English, we have used various resources, both textual and lexical. In the following sections we will briefly present these resources, what methodological framework was used for their development, and how they were integrated for their successful usage.

2.1 Textual Resources – Aligned Texts

The aligned texts as a special form of multilingual corpora were the focus of many projects in the past few decades. A systematic approach to the development of multilingual corpora was initiated within the Multext project, which subsequently included East-European languages through the Multext-East project [5]. In the meantime, many multilingual corpora have been compiled from large corpora, usually fully automatically prepared, which have a range comprised from texts in the limited technical domain [18] to more versatile literary corpora [5] that are often more modest in size but minutely prepared.

The main textual resource used to explore WS4LR is Jules Verne's novel *Around the World in Eighty Days*. This text was chosen for various reasons. First of all, the text is available in digital form for the majority of European languages, including Balkan languages. Regarding its content, it represents a suitable text for different types of analysis, especially in the domain of named entity recognition (geographical concepts and different measures). Besides this, it has already been used for some interesting research, e.g. multi-word tagging [13] and building models for machine translation [21]. Finally, from a practical point of view, its suitability stems from the fact that it presents a sample text for the French distribution of the Unitex system [15].

Versions of the novel in fifteen different languages have been acquired, but not all of these texts have yet been aligned. Among the already aligned texts are the French original and translations in English and four Balkan languages – Serbian, Bulgarian, Greek, and Romanian.

In the preparatory phase each translation was marked in accordance with the TEI-standard in XML, and the title (<head>), paragraph (<p>) and segments (<seg>) were included as units of a text logical layout. At the beginning of the alignment process, all segments coincided with sentences automatically tagged by Unitex. The XAlign system [1] was used for the alignment process. Starting from the French version, the goal of the alignment was to establish 1:1 relations with all other languages on the segmental level. In order to achieve this goal and after manually checking all aligned segments, some of them had to be divided into smaller units, and some were grouped into larger units. Thus, we arrived at the total of 4,409 segments in all texts. This way, the missing segments or the inconsistencies between the source text and its translations were identified in most of the cases. In the following example the English segment is given only for the sake of translation.

```
<tu id=" n2941">
  <seg lang="en">
    <s id="Verne80days.n2941">
      Between Omaha and the Pacific, the railway crosses a
      territory which is still infested by Indians and wild beasts, and a
      large tract which the Mormons, after they were driven from
      Illinois in 1845, began to colonise.</s></seg>
    <seg lang="fr">
      <s id="Verne80days.n2941">
        Entre Omaha et le Pacifique, le chemin de fer franchit une
        contrée encore fréquentée par les Indiens et les fauves, -- vaste
        étendue de territoire que les Mormons commencèrent à coloniser
        vers 1845, après qu'ils eurent été chassés de l'Illinois.</s></seg>
    <seg lang="sr">
      <s id="Verne80days. n2941">
        Između Omaha i Tihog okeana pruga prolazi kroz predeo
        u kome još ima Indijanaca i divljih zveri - prostranu zemlju koju
        su počeli naseljavati mormoni oko 1845. godine, kada su ih
        prognali iz države Illinois.</s> </seg>
    <seg lang="bg">
      <s id="Verne80days. n2941">
        Между Омаха и Тихия океан железопътната линия
        прекосява район, все още населяван от индианци и диви
        зверове. Това е обширна територия, която мормоните са
        започнали да колонизират около 1845 г., след като са били
        прогонени от щата Илинойс.</s></seg>
    <seg lang="gr">
      <s id="Verne80days. n2941">
        Ανάμεσα στην Ομάχα και στον Ειρηνικό, το τρένο
        διασχίζει περιοχές όπου συννάζουν ακόμα Ινδιάνοι και αγρίμια -
        τεράστια εδαφική έκταση την οποία αρχισαν να αποικίζουν οι
        μορμόνοι μετά το 1845, οπότε κυνηγήθηκαν από το
        Ιλινόις.</s></seg>
    <seg lang="ro">
      <s id=" Verne80days.n569">
```

```
    între Omaha și Pacific drumul de fier trece printr-o
    regiune populată încă de indieni și fiare, - vastă întindere pe care
    mormonii au început s-o colonizeze pe la 1845 după ce au fost
    izgoniți din Illinois.</s>
  </tu>
```

2.2 Morphological Dictionaries in LADL

Format

Morphological dictionaries are a necessary resource in various phases of the automatic analysis of a text. The tool WS4LR expects morphological dictionaries to be in the format known as DELAS/DELAF presented in [2] that was developed in LADL (*Laboratoire d'Automatique Documentaire et Linguistique*) under the guidance of Maurice Gross. The format of a DELAS-type dictionary basically consists of simple word lemmas accompanied with inflectional class codes which allow for the production of a DELAF-type dictionary, which consists of all inflectional forms with their grammatical information. In the Unitex environment, one finite-state transducer responsible for the generation of all inflectional forms of each DELAS lemma corresponds to each inflectional class code. The Serbian morphological dictionary of simple words contains 121,000 lemmas which yield the production of approximately 1,450,000 different lexical words. Close to 87,000 simple lemmas belong to the general lexica, while the remaining 34,000 lemmas represent various kinds of simple proper names [11]. The Bulgarian Grammar dictionary (DELAS dictionary) consists of 127,000 lemmas distributed as follows: app. 85,000 simple lemmas belong to the general lexis, app. 6,000 lemmas represent domain specific lexis and app. 36,000 lemmas are simple proper names. The corresponding DELAF dictionary consists of app. 1,260,000 entries [7].

2.3 Semantic Networks - Wordnet

Semantic networks, seen as one important node in the hierarchy of ontologies, are used more and more in various phases of the automatic analysis of texts. The tool WS4LR expects them to be in the form of wordnets, that is, nodes representing sets of synonymous words (synsets) which are linked by various semantic relations. The first built wordnet was an English wordnet, the so-called Princeton Wordnet (PWN), having today approximately 140,000 synsets. Due to its remarkable size and successful inclusion in various computer-based applications, it is considered as the de facto standard upon which wordnets for many other languages have been built. One successful application of this concept was achieved by the Balkanet project which was funded by the European Commission from (2001-2004). In the scope of this project, the development of wordnets for the following Balkan languages was initiated [20]: Bulgarian, Greek, Romanian, Serbian, and Turkish. The important feature of these wordnets is that they are all aligned with PWN via the Interlingual index (ILI) [22]. Namely, ILI consists of concepts, while wordnets represent

the lexicalization of concepts in various languages and the way which they are connected.

The Serbian wordnet today consists of more than 15,000 synsets built by app. 25,000 literals. All of them are linked to PWN, except for 532 Balkan specific concepts that are connected with other Balkan languages, and 155 Serbian specific concepts that remain unconnected with other languages. The Bulgarian wordnet consists of more than 31,000 synsets built by more than 66,000 literals. The synsets are linked with the PWN as well; again there are 436 Balkan specific concepts shared with other Balkan languages and 182 Bulgarian language specific concepts. Both Serbian and Bulgarian wordnets, as well as wordnets for other Balkan languages, are represented in WS4LR using common XML schema.

2.4 The Prolex Database

The *Prolex project* was initiated in the 1990's with the study of toponyms in French and had the aim of appropriately processing proper names in natural language applications [16]. This work was followed by the development of a Serbian version, which finally led to the design and construction of a relational multilingual dictionary of Proper Names, the Prolexbase, in the form of a relational database [19]. This model is based on two main concepts: the *pivot* (that represents the *conceptual proper name*) at a language independent level and the *prolexeme* (the projection of the pivot onto a particular language) which is a set of lemmas that includes the name, but also its aliases (variations in orthography, abbreviated forms, acronyms, etc.) and its derivatives. For instance, if a meronymy relation is established between the concepts "New York" and "the United States of America", then their Serbian Latin equivalents *Njujork* and *Sjedinjene Američke Države*, their Serbian Cyrillic equivalents *Њујорк* and *Сједињене Америчке Државе*, and their Bulgarian equivalents *Ню Йорк* and *Съединени американски щати* are all connected automatically.

3. Using WS4LR with Aligned Texts

The WS4LR module that works with aligned texts expects them to be in the Translation Memory eXchange (TMX) format¹. It can transform texts previously aligned by XAlign into this format as well as into several other formats: textual, XML and tabular. This is particularly important since XAlign has been integrated into Unitex software starting from its version 2.1. In addition, the user can also produce various visualizations of aligned texts by applying appropriate XSLT transformations. Thus, the user can freely browse with such visualized texts. One such visualization is represented in Figure 1.

¹ For details on TMX format see <http://www.lisa.org/tmx/tmx.htm>

Browsing, however, is not a particularly successful form of text exploration. The WS4LR module for aligned texts offers the user the ability to posit different forms of queries that can be automatically expanded by using various bilingual lexical resources presented in the previous section. WS4LR offers the user the possibility to expand the query not only morphologically and semantically, but also to another language. If the first language is Serbian, the second language can be English, Bulgarian, or any other. A user can choose two working languages by adjusting the parameters in the "Preferences" menu of WS4LR. Besides this, WS4LR provides further possibilities for the user to control their query formulation, since in addition to expansion it also offers the query to be narrowed. Namely, a user can reject some of the automatically offered query expansions.

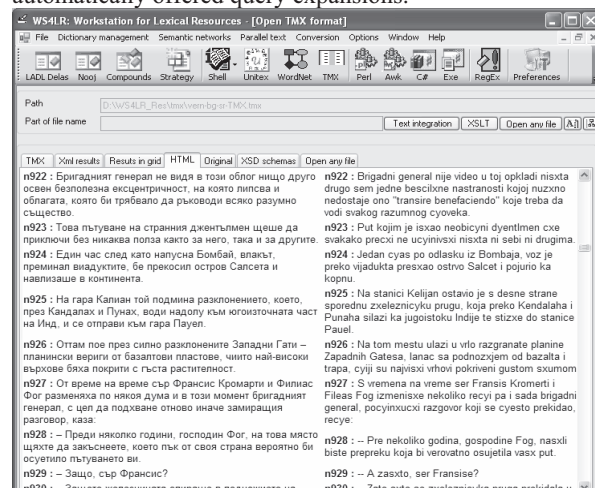


Figure 1. The HTML View of the Aligned Bulgarian-Serbian Text

User queries can be semantically expanded by the wordnets and by the Prolex database. WS4LR obtains the semantic expansion of a query by means of the wordnet of the first language (the Serbian wordnet – SWN, as is the case in our examples), selecting all synsets containing a given word and offering them to the user. This provides the user with an insight into all the concepts the keyword pertains to, through sets of synonyms used for these concepts. A user then gains the possibility to delete some of these synsets if he or she decides that they pertain to concepts which are not of interest at that particular moment. Also, the user can formulate a bilingual query by adding a second language to it. Namely, WS4LR can, for a given set of concepts, identify all corresponding concepts in the second language wordnet by using ILI. Thus, for an expanded Serbian query, one could obtain the corresponding expanded query in Bulgarian. The form used to bilingually expand a simple query *glava* "head" with the Bulgarian *глава* is presented in Figure 2. The semantic expansion is obtained by

checking the box “Semantic extension” in this form and by choosing the appropriate resource (Wordnet in this case), while the bilingual expansion is obtained by checking the box “Another language extension”.

In the same form, the user can choose to morphologically inflect all chosen keywords in both languages. If he or she wishes to do so the box “With inflection” should be checked. Morphological expansion is performed by Unitex modules that use morphological dictionaries of simple words as well as inflectional transducers. This option works only if a particular query keyword is listed in the morphological dictionary of the corresponding language. If this is not so, the aligned text will be searched only with the original keyword. As shown in Figure 2, the automatically added inflected forms of chosen keywords are presented in an editable form in which some of these inflected forms can be deleted or modified. For instance, the Serbian word *put* “path” has two plural forms: *putevi* and *puti*. The second one is restricted to poetical usage and a user can choose to delete it from the expanded query if the working text is not of that kind.

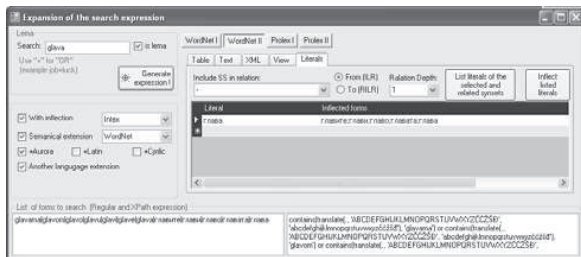


Figure 2. The original query keyword *glava* is shown in the upper left corner. The chosen query expansions are shown on the left side. The query expanded by the Bulgarian wordnet is shown on the right side, together with the automatically obtained list of inflected forms that can be edited. The two fields at bottom show the final query set.

Finally, when a query is launched, the result is obtained with all the retrieved occurrences highlighted (see Figure 3).

The query can be further semantically expanded by the choice of a particular semantic relation (e.g. hypernymy/hyponymy), in which case synsets pertaining to hypernyms/hyponyms of concepts from the initial group will also appear among the query set. This feature will be illustrated by a query which starts with the Serbian keyword *brodić* “small boat”. We would like to perform a bilingual search with a semantic expansion. The chosen Serbian keyword belongs to only one synset {brodica:1, brodić:1} whose corresponding Bulgarian synset is {лодка:1, ладия:1}. Figure 4 shows that these synsets are deep in the hypernymy/hyponymy hierarchy. In such a situation, expanding query with hypernym synsets can be useful.

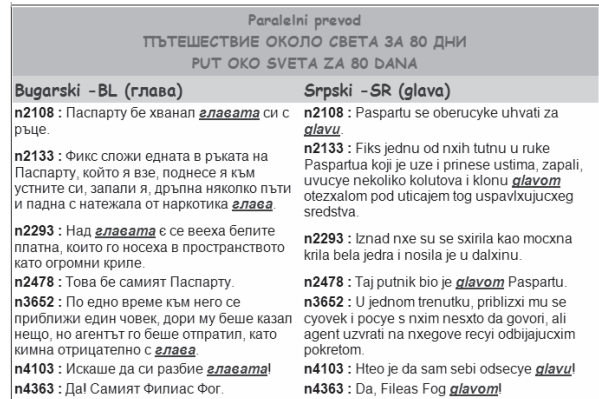


Figure 3. Some representative examples of aligned segments with the keywords *glava* and *glava* and their inflectional forms in HTML format.

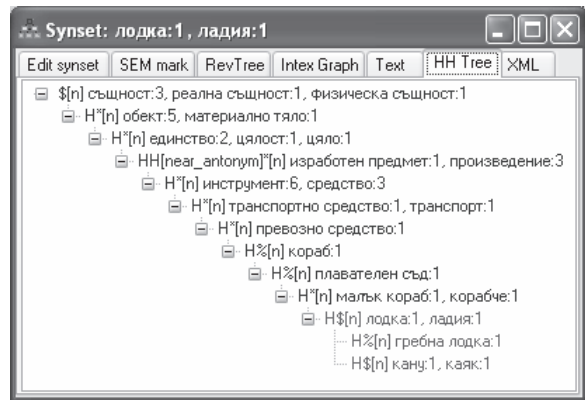


Figure 4. The hypernym/hyponym wordnet hierarchy of the Bulgarian synset {лодка:1, ладия:1}. The corresponding Serbian synset belongs to a similar tree.

Figure 5. shows the query expansion form in which the original query *brodić* is expanded, not only with a literal from its corresponding synset *brodica*, but also with the literals from synsets belonging to the hypernym branch of the length two, that are {barka:1, čamac:1, čun:1} “boat” and {lada:1} “vessel”.

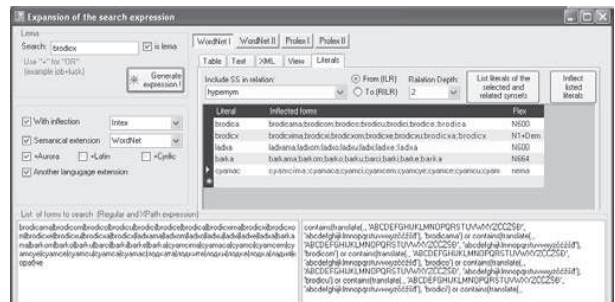


Figure 5. In the query expansion form, the user can choose the type of semantic relation for the expansion and the length of the path with this relation he or she wishes to pursue.

Since in this case a bilingual search is initiated, the user can perform the same semantic expansion for the second language, presented in Figure 6. The two Bulgarian literals thus obtained are *платателен съд* and *малък кораб* which are multi-word units. Since inflection of multi-word units for Bulgarian is not yet integrated in WS4LR, as will be explained in the final section, the user can choose to delete it from the final query set or to keep only the nouns *съд* and *кораб*, as we have done in our example search.

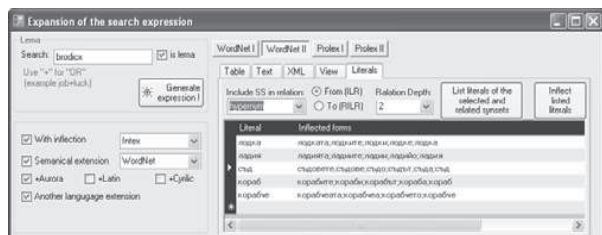


Figure 6. The semantic expansion in the second language – Bulgarian – using hypernym relation

The results obtained by this query are very interesting and, show the potential this tool offers for various linguistic and literary researches on their own. This query retrieved 129 aligned segments, each of which contained at least one of the keywords from the produced query set in at least one of the languages. It comes as a surprise that only 8 of these segments contained query keywords in both languages. This is mainly due to the fact that adjectives *платателен* and *малък* were omitted from the Bulgarian keywords thus broadening the query on the Bulgarian side too much. There were 5 segments with the keyword *съд*, with two occurrences of *платателен съд* “vessel”; none of them corresponded to the Serbian wordnet equivalent *lađa*. There were also 90 occurrences of *кораб* among which there was not one *малък кораб*; in this case, however, the Serbian equivalent for *кораб* was almost unmistakably *brod*, as suggested by both wordnets.

Figure 7 shows some examples of a partial retrieval. The first (n1616) and third (n2286) segments in this sample occur due to the fact that the reference to a “boat” is missing in one of the languages. The other segments show that the Serbian *brod*, besides corresponding to the English *ship* and the Bulgarian *кораб*, is also a generic notion and should probably be added to the hypernym synset (segments n2274, n2356 and n2439). On the other hand, Serbian *jedrilica* and *jedrenjak* “sailing vessel” are translated in Bulgarian with the “sister” synsets *кораб* or *корабче* instead of using a more specific Bulgarian word *платноход* (segments n2299 and n2323). In the last example (n3707), in Bulgarian a rather arbitrary choice *лодка* is made for a more specific type of a vessel referred to in Serbian as *kuter* “cutter”.

Paralelni prevod ПЪТЕШЕСТВИЕ ОКОЛО СВЕТА ЗА 80 ДНИ PUT OKO SVETA ZA 80 DANA	
Bugarski -BL (корабчето)	Srpski -SR (brodicx)
n1616 : Фикс го видя да слиза от файтона и да се качва в <i>лодка</i> с госпожа Ауда и прислужника му.	n1616 : Fiks ga vide gde silazi sa kola i gde se ukrcava sa gospodkom Audom i svojim slugom.
n2274 : Англиското знаме се вееше на нос на <i>корабчето</i> .	n2274 : Engleska zastava se leprskala na kixuni broda.
n2286 : – Лоцмане, мисля, че не е необходимо – каза Филмас Фог, когато излязоха в открито море – да ви напомням да се движите с възможно най-висока скорост.	n2286 : – Pilote, -- recyе Fileas Fog u trenutku kada je <i>brodicx</i> zaplovio na otvoreno more -- nije potrebno da vam preporucujem sxto vecxu zxurbu.
n2299 : Често срещаха кораби и със скоростта, с която се движеше, <i>лодкама</i> щеше да се разбие и при най-лекия удар.	n2299 : Sudari brodova tu nisu bili retki, a kako je jedrenjak jedrio velikom brzinom, razmrskao bi se pri najmanxem sudaru.
n2323 : На следващата заран, на 8 ноември, <i>корабчето</i> бе изминало над сто мили.	n2323 : Sutradan, 8. novembra, pri izlasku sunca jedrilica je bila preslxla visxe od stotinu milxa.
n2356 : По-сериозно не биха управлявали и <i>лодка</i> от “Кралския яхтклуб”.	n2356 : Ni na utakmici Kraljevskog jacht-kluba ne bi tacynije upravlxao brodom.
n2439 : Скоростта на <i>корабчето</i> бе забележителна.	n2439 : Brod je plovio dosta brzo.
n3707 : Те имат много платна, повече дори от състезателна <i>лодка</i> , и при попътен вятър се носат по заснежени полета с бързина, ако не по-голяма, то поне равна на скоростта на експресите.	n3707 : Uostalom, ona su odlicyno opremlxena jednima - bolxe nego neki gusarski kuter -, i kad im vetar duva s ledxa, klize po povrxsnini prejxe ako ne по-гольма, то поне равна на скоростта jednakoм brzinom mozxda i vecxom od brzih vozova.

Figure 7. A few examples of a partial retrieval

Figure 8 shows eight examples of the full retrieval. In one of these examples (n1972) for the Serbian *čamac* the near synonym in Bulgarian *корабчето* is used (as determined by wordnets). In two cases (n2267 and n2294) for the Serbian *brodic* the near hypernym *корабчето* is used, while in five cases (n514, n518, n586, n3827, n4049) for the Serbian *čamac* and *barka* the near hyponym *лодка* is used. This is not an unexpected result; rather it only proves that searching with the help of semantic networks, on the web for instance, can be useful, which is the ultimate goal of our experiments.

Bugarski -BL (корабчето)	Srpski -SR (brodicx)
n514 : Десетина <i>лодки</i> се отделиха от брега и се отправиха към “Монголия”.	n514 : Desetak <i>cyamacx</i> се otisnu od obale i uputi pred Mongoliju
n518 : По-голямата част обаче се качиха в <i>лодките</i> , които бяха приближили “Монголия”.	n518 : ali većxina се iskrcxa u <i>cyamcima</i> koji su pristali uz Mongoliju.
n586 : После се качи в една <i>лодка</i> , върна се на борда на “Монголия” и влезе в каютата си.	n586 : zatim се ukrcxa u jedan <i>cyamac</i> i vrati na Mongoliju, gde udxе u svoju kabinu.
n1972 : ... военни или търговски кораби; японски или китайски малки <i>корабчета</i> ; крайбрежни <i>корабчета</i> ...	n1972 : ... ratni trgovacycki brodovi; japanski i kineski <i>cyamci</i> , ...
n2267 : От задната част на <i>корабчето</i> се слизаше в квадратна каюта, ...	n2267 : Ispod krova straxnxeg dela <i>brodicx</i> silazilo се u cyetvrtastu odaju ...
n2294 : <i>Корабчето</i> , понасяно от вятъра, сякаш летеше във въздуха.	n2294 : <i>Brodicx</i> , nosxen vetrom, kao da je leteo vazduhom.
n3827 : Филмас Фог повика една <i>лодка</i> , качи се в нея и след няколко загребвания с веслата се озова пред стълбата на “Хенриета” – параход със железен корпус, чиято горна част беше дървена.	n3827 : Fileas Fog zakupi <i>cyamac</i> , sede u nhexa i posle nekoliko zaveslxaja nadxe се na lestvama Henrijeta, broda sa gvozdenim trupom kome je krov bio od drveta.
n4049 : “Хенриета” вече беше оголен и приличаше на понтонна <i>лодка</i> .	n4049 : Henrijeta je sad izgledala kao kakva pontonska <i>barka</i> .

Figure 8. All occurrences of a full retrieval

When a search is performed not with common keywords but with proper nouns then a query expansion with Prolex database offers more possibilities. Semantic relations incorporated into this database are adapted to proper names. Here, the user can choose to expand his query both on the conceptual and the linguistic level. It can be seen in Figure 9 how a query launched with a pivot *Paris* is linguistically expanded into two languages. A morphological expansion can be chosen here as well and it

is performed in the same way, using the same methods as for common words. In the given example, the query expansion for Serbian gives more results since the Prolex database for Bulgarian has only some sample entries.

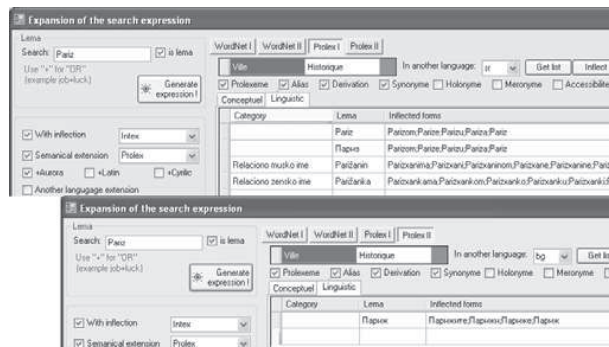


Figure 9. Prolex based semantic expansions

4. Additional Possibilities

We have illustrated the functions of WS4LR for working with aligned texts in the previous section by using the Serbian and Bulgarian pair. This can be successfully used for other Balkan languages as well. Wordnets have been developed through the Balkanet project for Greek, Romanian and Turkish, which have enabled the experiments to have semantic query expansions for those languages as well. For Greek [12] and Romanian [3], morphological dictionaries in the LADL format have also been developed – however, these resources were not at our disposal so we have not been able to experiment with morphological expansion for these languages.

The possibility and the need for some of the functions developed within WS4LR to become available also on the web have led to the development of the WS4QE web application for lexical resources. This application is still under development, but some of its functions can already be used. Numerous user functions are envisaged for this tool, but the largest set is related to the expansion of queries submitted to the Google search engine, and they have already been implemented. In fact, they are very similar to those presented in the previous section. The only difference is that expanded queries are not applied to an aligned text but are rather forwarded to the search engine.

Figure 10 shows such a retrieval that starts with the Serbian keyword barka “boat” and is further expanded by the Serbian synset {barka:1, čamac:1, čun:1} and Greek corresponding synset {βάρκα:0, λέμβος:0}. Figure 11 represents the first results retrieved with such an expanded query by Google.

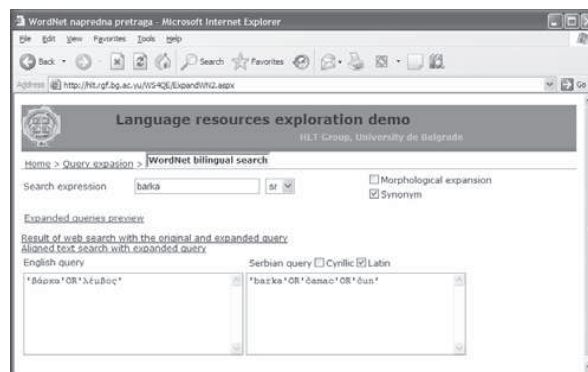


Figure 10. A bilingual query expansion with WS4QE – An example of Serbian and Greek

5. Further Work

Our main concern for our future work is the adequate processing of multi-word units. That is, we would like our tool to treat multi-word units in the same way as simple words and to inflect them correctly upon request. The first version of this approach was presented in [10]. Although this version gave promising results for Serbian, it was hardwired into the tool itself so that it was not easy to modify the Serbian module or to apply it to other languages. With a new approach that relies on the feature structure description of a particular language’s morphology [6] and widely uses XML technology, the portability to other languages will be much easier [17]. On a more practical level, our aim is enrich our lexical resources, first of all to enrich the Prolex database, as we plan to use it in a translation environment [14]. It is our wish to work in a future with a true aligned Balkan text – that is, a text originally written in a Balkan language and translated into other Balkan languages.

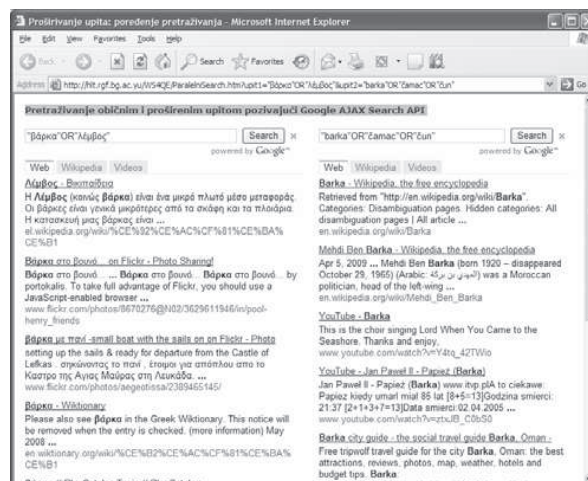


Figure 11. Results of a query bilingually expanded by Wordnet

6. References

- [1] P. Bonhomme, T. M. H. Nguyen, S. O'Rourke. XAlign: l'aligneur de Langue & Dialogue, <http://www.loria.fr/equipes/led/outils/ALIGN/align.html>, 2001.
- [2] B. Courtois, M. Silberstein (eds.). *Dictionnaires électroniques du français*. Langue française. 87, Larousse, Paris, 1990.
- [3] D.-M. Dimitriu. *Grammaires de flexion du roumain en format DELA*, Rapport interne 2005-02 de l'Institut Gaspard-Monge, CNRS, 2005.
- [4] T. Erjavec and N. Ide. The MULTEXT-East Corpus. In *LREC'98*, Granada, pp. 971-974, 1998.
- [5] A. Gelbukh, G. Sidorov, J.-A. Vera-Félix. A Bilingual Corpus of Novels Aligned at Paragraph Level. In proc. *FinTAL-2006. Lecture Notes in Artificial Intelligence*, no. 4139, Springer-Verlag, pp. 16–23, 2006.
- [6] ISO 24610. *Language resource management – Feature Structures*, ISO/TC 37/SC 4, 2005.
- [7] S. Koeva. *Modern language technologies – applications and perspectives*, in: *Lows of/for language*, Hejzal, Sofia, 2004, 111- 157, 2004.
- [8] C. Krstev, et al. Combining Heterogeneous Lexical Resources, in Proc. of the Fourth International Conference LREC, Lisbon, Portugal, May 2004, vol. 4, pp. 1103-1106, 2004.
- [9] C. Krstev, R. Stanković, D. Vitas, I. Obradović. *WS4LR: A Workstation for Lexical Resources*, Proceedings of the 5th International Conference on Language Resources and Evaluation, LREC 2006, Genoa, Italy, May 2006, pp. 1692-1697, 2006.
- [10] C. Krstev, R. Stanković, D. Vitas, I. Obradović, The Usage of Various Lexical Resources and Tools to Improve the Performance of Web Search Engines, in *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco, 28-30 May 2008, European Language Resources Association (ELRA), 2008.
- [11] C. Krstev. *Processing of Serbian*, Faculty of Phylology, University of Belgrade, Belgrade, 2008.
- [12] T. Kyriacopoulou. Les dictionnaires électroniques: Morphologie et syntaxe. Le cas du grec moderne, *Proceedings AILA 1990*, Chalcidique, 1990.
- [13] E. Laporte, T. Nakamura, S. Voyatzi. A French Corpus Annotated for Multiword Nouns, in: *Towards a Shared Task for Multiword Expressions (MWE 2008)*, in scope of the *Sixth International Conference on Language Resources and Evaluation (LREC'08)*, http://multiword.sourceforge.net/download/MWE2008-papers/8_Laporte.pdf, 2008.
- [14] D. Maurel, D. Vitas, C. Krstev, S. Koeva. Prolex: a lexical model for translation of proper names. Application to French, Serbian and Bulgarian, in *Bulag - Bulletin de Linguistique Appliquée et Générale*, Les langues slaves et le français : approches formelles dans les études contrastives, eds. A. Dziadkiewicz & I. Thomas, No. 32, pp. 55-72, Presses Universitaires de Franche Comté, Besançon, 2007.
- [15] S. Paumier. *Unitex 2.1 User Manual*, <http://www-igm.univ-mlv.fr/~unitex/UnitexManual2.1.pdf>, 2008.
- [16] O. Piton, D. Maurel. Beijing frowns and Washington takes notice: Computer Processing of Relations between Geographical Proper Names in Foreign Affairs, *Fourth International Workshop on Applications of Natural Language to Data Bases (NLDB'00)*, Versailles, 28-30 juin (Actes p. 66-78), 2000.
- [17] R. Stanković. Improvement of Queries using a Rule Based Procedure for Inflection of Compounds and Phrases. *Polibits (37) 2008, Special section: Natural Language Processing, Journal of Research and Development in Computer Science and Engineering*, ed. Grigori Sidorov, Centro Innovación y Desarrollo Tecnológico en Computo, Instituto Politécnico Nacional, Mexico, pp. 14-20, 2008.
- [18] R. Steinberger, B. Pouliquen, A. Widiger, C. Ignat, T. Erjavec, D. Tufiş. 2006. The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. In *Proceedings of the 5th LREC Conference*, Genoa, Italy, 22-28 May, 2006, pp.2142-2147, 2006.
- [19] M. Tran, D. Maurel. Prolexbase : Un dictionnaire relationnel multilingue de noms propres, *Traitement automatique des langues*, Vol. 47-3, 2006.
- [20] D. Tufiş (ed.). *Special Issue on BalkaNet Project*, Romanian Journal on Information Science and Technology. Bucureşti: Publishing house of the Romanian academy, Vol. 7, No.1-2, 2004.
- [21] D. Tufiş, S. Koeva, T. Erjavec, M. Gavrilidou, and C. Krstev. Building Language Resources and Translation Models for Machine Translation focused on South Slavic and Balkan Languages. In M. Tadić, M. Dimitrova-Vulchanova and S. Koeva (eds.) *Proceedings of the Sixth International Conference Formal Approaches to South Slavic and Balkan Languages (FASSBL 2008)*, pp. 145-152, Dubrovnik, Croatia, September 25-28, 2008.
- [22] P. Vossen (ed.) *EuroWordNet: A Multilingual Database with Lexical Semantic Networks*. Dordrecht: Kluwer Academic Publishers, 1998.