



**HAL**  
open science

## MULTEXT-East Resources for Serbian

Cvetana Krstev, Duško Vitas, Tomaž Erjavec

► **To cite this version:**

Cvetana Krstev, Duško Vitas, Tomaž Erjavec. MULTEXT-East Resources for Serbian. Zbornik 7. mednarodne multikonference Informacijska družba IS 2004 Jezikovne tehnologije 9-15 Oktober 2004, Ljubljana, Slovenija, 2004, Oct 2004, Ljubljana, Slovenia. hal-01108226

**HAL Id: hal-01108226**

**<https://hal.science/hal-01108226>**

Submitted on 22 Jan 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# MULTEXT-East Resources for Serbian

Cvetana Krstev,\* Duško Vitas,† Tomaž Erjavec‡

\*Faculty of Philology  
University of Belgrade  
Studentski trg 3, 11000 Begrade  
Serbia and Montenegro  
cvetana@matf.bg.ac.yu

†Faculty of Mathematics  
University of Belgrade  
Studentski trg 16, 11000 Begrade  
Serbia and Montenegro  
vitas@matf.bg.ac.yu

‡Department of Knowledge Technologies  
Jožef Stefan Institute  
Jamova 38, 1000 Ljubljana  
Slovenia  
tomaz.erjavec@ijs.si

## Abstract

MULTEXT-East is a multilingual dataset for language engineering research and development. This standardised and linked set of resources covers a large number of mainly Central and Eastern European languages and includes the EAGLES-based morphosyntactic specifications, defining the features that describe word-level syntactic annotations; medium scale morphosyntactic lexica; and annotated parallel, comparable, and speech corpora. The most important component is the linguistically annotated corpus consisting of Orwell's novel "1984" in the English original and translations. MULTEXT-East has already seen several editions, with the latest one being Version 3, where the most important addition has been that of Serbian language resources. The paper presents MULTEXT-East Version 3 with special emphasis on the Serbian components, namely the structurally annotated "1984", the morphosyntactic specifications, the morphosyntactic lexicon and the linguistically annotated "1984". The complete dataset, unique in terms of languages and the wealth of encoding, is extensively documented, and freely available for research purposes.

## Jezikovni viri MULTEXT-East za srbski jezik

MULTEXT-East je večjezikovna podatkovna množica, namenjena raziskavam in razvoju jezikovnih tehnologij. Ta standardizirana in povezana množica jezikovnih virov pokriva veliko število predvsem srednje- in vzhodnoevropskih jezikov in vsebuje (1) oblikoslovne specifikacije, ki definirajo oznake za opis skladenjskih lastnosti besed, (2) oblikoslovne leksikone srednje velikosti in (3) označene vzporedne, primerljive in govorjene korpuse. Najpomembnejša komponenta je jezikovno označen korpus, ki vsebuje roman "1984" G. Orwella v angleškem originalu in prevodih. MULTEXT-East je doživel že več izdaj, pri čemer je zadnja t.i. verzija 3, kjer so glavni dodatek viri za srbski jezik. Članek predstavi MULTEXT-East verzijo 3 s posebnim poudarkom na srbskih komponentah, in sicer na strukturno označenem besedilu romana "1984", oblikoslovnih specifikacijah, oblikoslovnem leksikonu in jezikovno označenem besedilu "1984". Celotna podatkovna množica, enkratna glede na vsebovane jezike in bogastvo oznak, je podrobno dokumentirana in prosto dostopna v raziskovalne namene.

## 1. Introduction

The mid-nineties saw – to a large extent via EU projects – the rapid development of multilingual language resources and standards for human language technologies (Armstrong et al., 1998; Ide and Véronis, 1994; EAGLES, 1996). However, while the development of resources, tools, and standards was well on its way for EU languages, there had been no comparable efforts for the languages of Central and Eastern Europe. The MULTEXT-East project (Multilingual Text Tools and Corpora for Eastern and Central European Languages) was a spin-off of the EU MULTEXT project (Ide and Véronis, 1994); MULTEXT-East ran from '95 to '97 and developed standardised language resources for six CEE languages (Dimitrova et al., 1998), as well as for English, the 'hub' language of the project. The

project also adapted existing tools and standards to these languages. The main results of the project were lexical resources and an annotated multilingual corpus. The most important resource turned out to be the parallel corpus — heavily annotated with structural and linguistic information — which consists of Orwell's novel "1984" in the English original and translations.

One of the objectives of MULTEXT-East has been to make its resources freely available for research purposes. In the scope of the TELRI concerted action (Trans European Language Resources Infrastructure), the results of MULTEXT-East had been extended with several new languages and first released on a CD-ROM, and later through Web download via TRACTOR, the TELRI Research Archive of Computational Tools and Resources.

The Serbian language did not have its representative in the MULTEXT-East project. The researchers from the Faculty of Mathematics, however, participated in the TELRI concerted action. One of the results of this participation was the Serbian “1984” structurally annotated corpus, but the morphosyntactic specification, lexicon and linguistically tagged “1984” were not produced.

Following the TELRI release, the MULTEXT-East resources were used in a number of studies and experiments. In the course of such work, errors and inconsistencies were discovered in the MULTEXT-East specifications and data, most of which were subsequently corrected. But because this work was done at different sites and in different manners, the encodings of the resources had begun to drift apart.

The ’98–’00 EU Copernicus project CONCEDE (Consortium for Central European Dictionary Encoding) offered the possibility to bring the versions back on a common footing. Although CONCEDE was primarily devoted to machine readable dictionaries and lexical databases, one of its workpackages did consider the integration of the dictionary data with the MULTEXT-East corpus (Erjavec et al., 2003a). The CONCEDE release contained the revised and expanded morphosyntactic specifications, the revised lexica, and the significantly corrected and re-encoded linguistically annotated “1984” corpus.

In addition to delivering resources per-se, a focus of the MULTEXT-East, TELRI and CONCEDE projects was also the adoption and promotion of encoding standardisation. On the one hand, the morpholexical annotations and lexica were developed in the formalism of the (EAGLES-based) specifications for six Western European languages of the MULTEXT project (Ide and Véronis, 1994). On the other, in the TELRI edition, all the corpus resources were encoded in SGML, in CES, the Corpus Encoding Standard (Ide, 1998). For the corpus taken forward into the second edition, the Text Encoding Initiative Guidelines were adopted, in particular TEI P3 (Sperberg-McQueen and Burnard, 1999).

Finally, in 2004 the third version of the MULTEXT-East resources was released (Erjavec, 2004). This release offers several contributions: it brings together the first two, i.e., offers both the TELRI and CONCEDE versions in one package; all the resources have been recoded in XML, according to TEI P4 (Sperberg-McQueen and Burnard, 2002), thus enabling them for processing with XML-based tools; and resources for new languages have been added, in particular the morphosyntactic specification for Resian, a dialect of Slovene, and, crucially the morphosyntactic specification and the annotated Orwell for Serbian.

Version 3 also contains extensive documentation, e.g., navigational HTML pages, which serve to structure and link the resources, and which include the list of participants and indexes to the resource by type and language. While the TEI headers give the most precise and up-to-date information on the corpus components, the documentation also contains a bibliography with copies of the MULTEXT-East project reports (giving details of the resources, e.g., the corpus markup process), published papers, a mirror of the TEI P4 and CES documentation and certain related MULTEXT and EAGLES reports.

A complete description of the Version 3 resources is

```
<text id="mteo-sr." lang="sr">
<body id="Osr" lang="sh">
<div id="Osr.1" n="1" type="part">
<head>Prvi deo</head>
<div id="Osr.1.2" n="1" type="chapter">
<head>1.</head>
<p id="Osr.1.2.2">
<s id="Osr.1.2.2.1">Bio je vedar i
hladan aprilski dan; na &#x10D;asovnicima
je izbijalo trinaest.</s>
<s id="Osr.1.2.2.2"><name>Vinston
Smit</name>, brade zabijene u nedra da
izbegne ljuti vetar, hitro zama&#x10D;e u
staklenu kapiju stambene zgrade
<hi rend="it">Pobeda</hi>, no nedovoljno
hitro da bi spre&#x10D;io jednu spiralu
o&#x161;tre pra&#x161;ine da u&#x111;e
zajedno s njim.</s>
</p>
...
```

Figure 1: The structurally annotated Orwell

given in (Erjavec, 2004) and in the on-line documentation, while this paper concentrates on the Serbian resources. In the next section we introduce the structurally annotated Serbian “1984” (already a part of the TELRI release), in Section 3 we describe INTEX, the system that has for a long time served as the infrastructure for developing LR resources for Serbian, Section 4 explains the MULTEXT-East (Serbian) morphosyntactic specification, Section 5 the linguistically annotated “1984”, Section 6 the Serbian lexicon and the last section gives some conclusions and direction for further work.

## 2. Structural “1984” and alignments

The MULTEXT-East multilingual parallel corpus consists of the novel “1984”, about 100,000 words in length. The corpus contains extensive headers and markup for document structure, sentences, and various sub-sentence annotations, which have been harmonised over languages. As an example, the start of the text from the Serbian Orwell is given in Figure 1.

The translations of “1984” have been automatically sentence aligned with the English original, and the alignments hand-validated. The bilingual alignments are stand-off, i.e., they are stored not with the primary data but in separate documents, as references to sentence IDs; a (hypothetical) example is given in Figure 2.

The cesDoc encoded novel served as the basis for producing the linguistically annotated version. The link between the two is maintained via sentence identifiers.

The Serbian version was produced already in the scope of TELRI. The digital source was the same as for the English and Slovene versions, namely the Oxford Text Archive, via the ECI multilingual CD-ROM. This version was plain ASCII, so it was first marked up, similar to other versions in SGML, and then sentence segmented and aligned with the English original. Also, many typographical errors were corrected.

```

<?xml version="1.0" encoding="us-ascii"?>
<!DOCTYPE cesAlign SYSTEM "xcesAlign.dtd">
<cesAlign version="4.1">
  <linkList id="Osren">
    <linkGrp id="Osren.1" type="body"
      targettype="s" domains="Osr Oen">
      <link xtargets="Osr.1.2.2.1 ;
        Oen.1.1.1.1/"> <!--1:1-->
      <link xtargets="Osr.1.2.2.2 ;
        Oen.1.1.1.2
        Oen.1.1.2.1/"> <!--1:2-->
      <link xtargets=";
        Oen.1.1.2.2/"> <!--0:1-->
    ...

```

Figure 2: Stand-off sentence alignments

### 3. INTEX and Serbian resources

Before discussing the MULTEXT-East Serbian morphosyntactic resources (the specifications, lexicon and linguistically annotated “1984”) we first describe the basis for these resources, which had been developed independently of European projects, namely the Serbian morphological lexicon (Vitas and Krstev, 2001) in the format of the INTEX system, which is based on the technology of finite-state transducers (Silberstein, 2000).

In this dictionary a lemma is of the form:

$$W_t, W_l.Cn + SSD : (Codes)*$$

where  $W_t$  represents the textual word,  $W_l$  the corresponding lexical word,  $C$  is the part of speech,  $n$  is the code of inflective class,  $SSD$  is the set of syntactic and semantic attributes of the lemma that is classified as  $Cn$ , and  $codes$  describe the values of morphological categories that realized with the form  $W_t$ . For instance, the dictionary entry of the noun *prozor* (Engl. window) is *prozorom, prozor.N01+Com:ms6q*, which describes the form *prozorom* as the form of common (+Com) masculine (m) inanimate (q) noun (N) from the inflective class 01 in instrumental (6) of singular (s). It can be seen that this format is not as compact as MSD, as the relevant information is distributed among the inflective code, syntactic and semantic information, both associated to lexical word, and the grammatical codes which are assigned to the textual word. The *codes* are not positional — for a part of speech one alphanumeric character represents a value of one and only one of its attributes.

The present size of the Serbian morphological dictionary, given in Table 1 enables morphological text analysis with a high percentage of success, around 90% for literary texts. Some of the word-forms that are not covered by the dictionary itself can be successfully morphologically tagged by additional tools (lexical transducers) incorporated in Intex. The use of this specifically constructed set of lexical transducers enables the recognition of various derived word-forms, such as several classes of compounds, possessive adjectives, diminutives, augmentatives, etc. (Pavlović-Lažetić et al., 2004).

The team from the University of Belgrade plans to convert its full INTEX lexicon to a MSD-type lexicon. It is also planned to tag with MSDs the corpus of contemporary

PoS	Lemmas	Word-Forms
Nouns	31,000	185,000
Verbs	14,660	436,000
Adjectives	21,800	352,000
other	3,400	7,600
Total	70,860	980,600

Table 1: The current size of the Serbian morphological dictionary

Serbian that is being developed at the Faculty of Mathematics (Vitas et al., 2003), where it plans to use INTEX and the lexica incorporated in it as a preprocessor.

### 4. Morphosyntactic Specifications

The MULTEXT-East morphosyntactic specifications give the syntax and semantics of the morphosyntactic descriptions (MSDs) used in the lexica and corpora. The MSDs, are structured and more detailed than is commonly the case for part-of-speech tags; they are compact string representations of a simplified kind of feature structures. The first letter of a MSD encodes the part of speech, e.g., Noun or Adjective. The letters following the PoS give the values of the position determined attributes. The specifications define, for each part of speech, its appropriate attributes, their values and one-letter codes. So, for example, the *Ncmpi* MSD expands to *PoS:Noun, Type:common, Gender:male, Number:plural, Case:instrumental*. It should be noted that in case a certain attribute is not appropriate (1) for a language, (2) for the particular combination of features, or (3) for the word in question, this is marked by a hyphen in the attribute’s position. Slovene verbs in the indicative, for example, are not marked for gender or voice, hence the two hyphens in *Vcip3s--n*.

The specifications have been developed in the formalism and on the basis of specifications for six Western European languages of the EU MULTEXT project (Ide and Véronis, 1994) and in cooperation with EAGLES, the Expert Advisory Group on Language Engineering Standards. Originally, these specifications were released as a report of the MULTEXT-East project but have been revised for both subsequent releases, and have become, if not a standard, then at least a reference for comparison (Erjavec et al., 2003b).

The MULTEXT-East morphosyntactic specifications have the following structure: (1) introductory matter; (2) the common specification; and (3) a language particular section for each language.

The common part of the specifications first defines the parts of speech and their codes; MULTEXT-East distinguishes the following, where not all PoS are used for all languages - we mark in *italics* those that are not used for Serbian: Noun (N), Verb (V), Adjective (A), Pronoun (P), *Determiner (D)*, *Article (T)*, Adverb (R), Adposition (S), Conjunction (C), Numeral (M), Interjection (I), *Residual (X)*, Abbreviation (Y), and Particle (Q).

The formal core of the specifications resides in the common tables; they define the features, their codes for MSD

```

<fLib type="Verb">
<f id="V0."
  select="en ro sl cs bg et hu hr sr sl-rozaj"
  name="PoS"><sym value="Verb"/></f>
<f id="V1.m"
  select="en ro sl cs bg et hu hr sr sl-rozaj"
  name="Type"><sym value="main"/></f>
<f id="V1.a"
  select="en ro sl cs bg et hu hr sr sl-rozaj"
  name="Type"><sym value="auxiliary"/></f>
<f id="V1.o"
  select="en ro sl cs et hr sr sl-rozaj"
  name="Type"><sym value="modal"/></f>
...

```

Figure 4: Morphosyntactic specifications as TEI features

representation, and their appropriateness for each language — an example is given in Figure 3.

Technically, the complete specifications are a  $\LaTeX$  document (with derived Postscript, PDF and HTML renderings), where the common tables are plain ASCII in a strictly defined format. This format is suitable for a printed version, tolerable for one in HTML, and reasonably manageable for modification and addition of new languages. However, it is not suitable for processing needs, in particular to enable smooth manipulation and linking to an XML encoded corpus using the MSDs.

We have therefore implemented a (Perl) conversion of the common tables into XML, using the TEI.fs module, a tagset devoted to encoding feature-structures. This tagset is currently being used as the basis of an evolving ISO standard (currently a Draft International Standard), as part of work of ISO/TC 37/SC4 Language Resource Management.

The XML version of the common tables has one feature library for each category, e.g.,  $\langle fLib\ type="Noun" \rangle$ . Each feature in such a library is comprised of the identifier, which enables the linkage to corpus MSDs, the name of the attribute, the languages the feature is appropriate for, and the symbol that is its value; examples are given in Figure 4.

The Serbian specifications was produced on the basis of the Croatian one (which was added in the scope of CONCEDE), with some modifications stemming less from the differences between languages and more by the set of morphosyntactic attributes already incorporated in the Intex e-dictionaries for Serbian. For instance, in the verb table the value gerund for the attribute VForm is the most appropriate to account for present and past gerund active in Serbian. Also, the attribute Clitic is applicable to Serbian copula verbs, as well as the attribute Aspect to the most of the other verbs. Both these attributes are already encoded for all verbs in Serbian e-dictionary. One of the other differences between Croatian and Serbian tables is the recognition of the value ‘paukal’ for the attribute ‘Number’ for several PoS in Serbian.

## 5. Lexicons

The MULTEXT-East morphosyntactic lexicons have a simple structure, where each lexical entry is composed of three fields:

1. the *word-form*, which is the inflected form of the word,

```

<fsLib type='Verb'>
<fs id="Van"
  select="en et"
  feats="V0. V1.a V2.n"/>
<fs id="Van----an----n"
  select="cs"
  feats="V0. V1.a V2.n V7.a V8.n V13.n"/>
<fs id="Van----an-n---p"
  select="sr"
  feats="V0. V1.a V2.n V7.a V8.n V10.n
        V14.p"/>
<fs id="Van----ay----n"
  select="cs"
  feats="V0. V1.a V2.n V7.a V8.y V13.n"/>
<fs id="Vanp"
  select="ro"
  feats="V0. V1.a V2.n V3.p"/>
...

```

Figure 5: MSDs as TEI feature structures

as it appears in the text, modulo sentence-initial capitalisation;

2. the *lemma*, which is the base-form of the word; where the entry is itself the base-form, the lemma is given as the equal sign; and
3. the *MSD*, i.e., the morphosyntactic description.

To produce the lexica, the token lists of the MULTEXT-East corpus were first fed through morphological analysers in order to produce the lemma list; this list was further extended from the comparable corpus, to arrive at at least 15,000 lemmas – some languages have further extended this, e.g., Romanian to 41,000 lemmas. In the next step, these lemmas were fed back to morphological generators (except for the agglutinative languages) in order to produce the complete inflected lists, i.e., the full paradigms of the lemmas, which constituted the final lexica of the project.

The MULTEXT-East lexica serve as medium sized morphological lexica for the languages. In addition to explicating the inflectional behaviour of the most common (and, typically, morphologically the most complex) words of the languages, the lexica also serve to establish the definitive set of valid MSDs for the languages.

For Serbian, currently, only a minimal lexicon was produced, which contains just the word-forms that in fact appear in the annotated “1984” corpus. This lexicon has 20,294 entries, 16,907 different word-forms, 8,392 lemmas and 906 MSDs.

To serve as a standard registry of MSDs, we converted the lexical MSDs to TEI feature structure libraries,  $\langle fsLib \rangle$ , one for each category. Here each MSD is expressed as a feature structure specifying its *id*, the language(s) it is appropriate for, and its decomposition into features. Some examples are given in Figure 5.

Both the  $\langle fsLib \rangle$ s and the  $\langle fLib \rangle$ s are stored in dedicated  $\langle TEI.2 \rangle$  element, complete with its TEI header. This document also constitutes a part of the linguistically annotated MULTEXT-East corpus.

### 3.2 Verb (V)

= =====			EN	RO	SL	CS	BG	ET	HU	HR	SR	SL-ROZAJ
P	ATT	VAL	C	x	x	x	x	x	x	x	x	x
= =====												
1	Type	main	m	x	x	x	x	x	x	x	x	x
		auxiliary	a	x	x	x	x	x	x	x	x	x
		modal	o	x	x	x		x		x	x	x
		copula	c		x	x	x			x	x	x
		base	b	x								
- - - - -			-									
2	VForm	indicative	i	x	x	x	x	x	x	x	x	x
		subjunctive	s		x							x
		imperative	m		x	x	x	x	x	x	x	x
		conditional	c	x		x	x	x	x	x	x	
		infinitive	n	x	x	x	x	x	x	x	x	x
		participle	p	x	x	x	x	x		x	x	x
		gerund	g		x			x	x		x	
		supine	u			x		x				x
		transgressive	t				x					
		quotative	q					x				
- - - - -			-									

Figure 3: Start of Common Table for Verbs

## 6. Linguistically annotated “1984”

The centrepiece of the MULTEXT-East resources is the linguistically annotated “1984”; it contains word level markup, namely context disambiguated lemmas and MSDs. Because it was the first such resources for many of the MULTEXT-East languages, also Serbian, it was the most difficult and time-consuming to produce as the work had to proceed mostly manually. The annotated novel is useful as a dataset for tagger and lemmatiser induction and testing, and has already been used for this purpose in a number of experiments: an evaluation exercise (Džeroski et al., 2000) compared four state-of-the-art trainable taggers on “1984”; (Hajič, 2000) tested a feature-based tagger on the corpus; and research on tagset reductions was investigated in developing tagging models for Romanian (Tufiş, 1999) and Hungarian (Varadi and Oravecz, 1999). Another strand of research used the corpus to investigate inductive learning of rules for morphological analysis, in order to lemmatise unknown words in a text (Erjavec and Džeroski, 2004). The annotated “1984” corpus also served as a testbed for experiments in word sense disambiguation (Ide et al., 2002) where the use of translation equivalents for automatic sense-tagging was investigated.

The work on the Serbian annotation proceeded in the following steps. First, using Intex as the tool and all of the Serbian lexical resources Serbian, “1984” was morphologically tagged. As a result, a textual file is obtained that contains the finite automaton of the text represented in the form of a regular expression. As an example, the first two sentences are given in Figure 6. In this file, each line contains one running text word, with zero, one or more associated lexical words. If no lexical word is associated this means that the word-form could not be recognized by the lexical resources. To each recognised word-form a set of grammatical categories is associated which represents the possible realizations of the word in text (for instance, for nouns: case, number, gender, and animacy).

```
{S}{Bio,biti.V77:Gsm}
({je, jesam.V575+Imperf+It+Iref+Aux:Pzsi}
 + {je,on.PRO+Prs:sz2fi:sz4fi})
{vedar,.A18:akms1g:akms4q}
({i,.CONJ} + {i,.PAR})
{hladan,.A18:akms1g:akms4q}
{aprilski,.A2+PosQ
 :adms1g:aems4q:aems5g:aemp1g:aemp5g}
({dan,.A1+PP:akms1g:aems4q}
 + {dan,dati.V103+Perf+Tr+Iref+Ref:Tms})
;
{S}
({na,.PREP+p4} + {na,.PREP+p7})
{cyasovnicima,.?}
({je, jesam.V575+Imperf+It+Iref+Aux:Pzsi}
 + {je,on.PRO+Prs:sz2fi:sz4fi})
{izbijalo,izbijati.V101+Perf+Tr+It+Iref:Gsn}
{trinaest,.NUM}
.
```

Figure 6: The output of the annotation by Intex

In the second step the text annotated in this way was manually checked and disambiguated. It means that not only have right lemmas and morphological categories been chosen for ambiguous word-forms and added for those words that had not been recognized, but non-ambiguous forms have also been checked in case that they had been incorrectly recognized (as is the case with the word-form *dan*, the last word of the first sentence in Figure 6). As a result, a non-ambiguous representation of a text is obtained in a same format. For instance, the regular expression for the same two sentences after manual correction and disambiguation is shown in Figure 7. This step had been done iteratively, chapter by chapter, which enabled both the correction of the used dictionaries and other lexical resources, and their enhancement.

In the third step, a Perl script was written and used to convert the Intex annotated text with to the MULTEXT-

```

{S}{Bio,biti.V77:Gsm}
{je,jesam.V575+Imperf+It+Iref+Aux:Pzsi}
{vedar,.A18:akms1g}
(i,.CONJ)
{hladan,.A18:akms1g}
{aprilski,.A2+PosQ:adms1g}
{dan,.N51:ms1q}
;
{S}
{na,.PREP+p7}
{cyasovnicima,cyasovnik.N9:mp7q}
{je,jesam.V575+Imperf+It+Iref+Aux:Pzsi}
{izbijalo,izbijati.V101+Perf+Tr+It+Iref:Gsn}
{trinaest,.NUM}
.

```

Figure 7: Textual file containing the disambiguated text automaton in Intex format

East annotation. The conversion is not, however, a straightforward task, not only because of the different encoding systems, as described in section 3., but also because of the differently chosen attributes. This difference can be most easily described in the case of verbs. For verbs, in MULTEXT-East the second attribute specifies a verb form, and the third a tense. However, due to the composite tenses, some verb forms are used for the construction of different tenses. For instance, in Serbian, verb form *imao* is the active past participle of the verb *imati* (Engl. *to have*), and is used to produce both perfect tense if used with the indicative form of the present tense of the copula verb *biti* (Engl. *to be*) — Figure 8 a) — and conditional if used with the conditional form of the same copula verb — Figure 8 b). In Intex, however, only the verb forms are recognized, which in the case of simple tenses enables the recognition of a tense as well (for instance, for present or aorist). For analytical tenses, the word-form recognition is not enough and more complex tools have to be used, as described in (Vitas and Krstev, 2003). These tools, as not yet being fully developed, were not used in the first step of the annotation process, and thus the precise mapping from Intex to MULTEXT-East tags was not possible. As a consequence, despite having different functions, the active past participle is always given the same value in the third attribute, that is `tense=past`, as being the most frequent.

The TEI P4 markup of the linguistically annotated Serbian “1984” obtained through this process is exemplified in Figure 9 by the same two sentences.

We also plan to use the annotated “1984” in the scope of the BalkaNet project for the validation of the Serbian WordNet being produced, along with the other languages involved in both MULTEXT-East and BalkaNet, i.e., Czech, Bulgarian and Romanian (Krstev et al., 2004).

## 7. Conclusions

The paper presented Version 3 of the MULTEXT-East resources, and, in particular, its Serbian language resources. As the resources cover a number of inflectionally rich languages, are interlinked, harmonised, have a standardised encoding, and have been manually validated and tested in practice, they can serve as a “gold standard” dataset for language technology research and development.

```

a) ...
<w lemma="on" ana="Pp3msn">on</w>
<w lemma="jesam"
ana="Va-p3s-an-y---p">je</w>
<w lemma="vecx" ana="Q-">vecx</w>
<w lemma="imati"
ana="Vmcs-sman-n---p">imao</w>
<w lemma="naslaga"
ana="Ncfsa--n">naslage</w>
<w lemma="salo" ana="Ncnsg--n">sala</w>
...
b) ...
<w lemma="ko" ana="C-s">Ko</w>
<w lemma="god" ana="Q-">god</w>
<w lemma="biti"
ana="Vmca3s-an-n---p">bi</w>
<w lemma="imati"
ana="Vmcs-sman-n---p">imao</w>
<w lemma="u" ana="Sps-">u</w>
<w lemma="ruka" ana="Ncfsl--n">ruci</w>
<w lemma="dokument"
ana="Ncmsa--n">dokument</w>
...

```

Figure 8: a) the active past participle *imao* used in active voice perfect tense; b) the active past participle *imao* used in conditional I

```

<text lang="sr" id="Osr.">
<body>
<div id="Osr.1" type="part" n="1">
<div id="Osr.1.2" type="chapter" n="1">
<p id="Osr.1.2.2">
<s id="Osr.1.2.2.1">
<w lemma="biti"
ana="Vmcs-sman-n---p">Bio</w>
<w lemma="jesam"
ana="Va-p3s-an-y---p">je</w>
<w lemma="vedar"
ana="Afpmsnn">vedar</w>
<w lemma="i"
ana="C-s">i</w>
<w lemma="hladan"
ana="Afpmsnn">hladan</w>
<w lemma="aprilski"
ana="Aopmpn">aprilski</w>
<w lemma="dan"
ana="Ncmsn--n">dan</w>
<c>i</c>
<w lemma="na"
ana="Spsa">na</w>
<w lemma="&#x10D;asovnik"
ana="Ncmsa--n">&#x10D;asovnicima</w>
<w lemma="jesam"
ana="Va-p3s-an-y---p">je</w>
<w lemma="izbijati"
ana="Vmcs-snan-n---e">izbijalo</w>
<w lemma="trinaest"
ana="Mc---l">trinaest</w>
<c>.</c>
</s>
...

```

Figure 9: The linguistic annotation of “1984”

While portions of the resources are distributed without any restrictions, the resources as a whole are available free of charge for research purposes only, as this was the condition imposed by some copyright holders of the sources.

Version 3 of the resources can be downloaded from the MULTEXT-East home page, <http://nl.ijs.si/ME/>. Access is enabled by filling out and submitting a Web based agreement, which is modelled after the one used by Edinburgh's Language Technology Group.

Currently, there are no plans to start working on Version 4; rather, the focus will be on the utility of V3, in our own research, and in enabling others to use the resources, by providing maintenance, continuing to support their accessibility and correcting errors.

### Acknowledgments

The work presented in this paper was, in part, supported by the bi-lateral project on scientific and technological cooperation between Slovenia and Serbia "The development of language resources for machine translation between the Slovene and Serbian languages".

### 8. References

- Armstrong, Susan, Masja Kempen, David McKelvie, Dominic Petitpierre, Reinhardt Rapp, and Henry Thompson, 1998. Multilingual corpora for cooperation. In *First International Conference on Language Resources and Evaluation, LREC'98*. Granada: ELRA.
- Dimitrova, Ludmila, Tomaž Erjavec, Nancy Ide, Heiki-Jan Kaalep, Vladimír Petkevič, and Dan Tufiş, 1998. Multext-East: Parallel and Comparable Corpora and Lexicons for Six Central and Eastern European Languages. In *COLING-ACL '98*. Montréal, Québec, Canada.
- Džeroski, Sašo, Tomaž Erjavec, and Jakub Zavrel, 2000. Morphosyntactic Tagging of Slovene: Evaluating PoS Taggers and Tagsets. In *Second International Conference on Language Resources and Evaluation, LREC'00*. Paris: ELRA. [Http://nl.ijs.si/et/Bib/LREC00/lrec-tag-www/](http://nl.ijs.si/et/Bib/LREC00/lrec-tag-www/).
- EAGLES, 1996. Expert advisory group on language engineering standards. [Http://www.ilc.pi.cnr.it/EAGLES/home.html](http://www.ilc.pi.cnr.it/EAGLES/home.html).
- Erjavec, Tomaž, 2004. MULTEXT-East Version 3: Multilingual Morphosyntactic Specifications, Lexicons and Corpora. In *Fourth International Conference on Language Resources and Evaluation, LREC'04*. Paris: ELRA. [Http://nl.ijs.si/et/Bib/LREC04/](http://nl.ijs.si/et/Bib/LREC04/).
- Erjavec, Tomaž and Sašo Džeroski, 2004. Machine Learning of Language Structure: Lemmatising Unknown Slovene Words. *Applied Artificial Intelligence*, 18(1):17–41.
- Erjavec, Tomaž, Roger Evans, Nancy Ide, and Adam Kilgarriff, 2003a. From Machine Readable Dictionaries to Lexical Databases: the Concede Experience. In *Proceedings of the 7th International Conference on Computational Lexicography, COMPLEX'03*. Budapest, Hungary.
- Erjavec, Tomaž, Cvetana Krstev, Vladimír Petkevič, Kiril Simov, Marko Tadić, and Duško Vitas, 2003b. The MULTEXT-East Morphosyntactic Specifications for Slavic Languages. In *Proceedings of the EACL 2003 Workshop on Morphological Processing of Slavic Languages*. Budapest.
- Hajič, Jan, 2000. Morphological Tagging: Data vs. Dictionaries. In *ANLP/NAACL 2000*. Seattle.
- Ide, Nancy, 1998. Corpus Encoding Standard: SGML guidelines for encoding linguistic corpora. In *First International Conference on Language Resources and Evaluation, LREC'98*. Granada: ELRA. [Http://www.cs.vassar.edu/CES/](http://www.cs.vassar.edu/CES/).
- Ide, Nancy, Tomaž Erjavec, and Dan Tufiş, 2002. Sense Discrimination with Parallel Corpora. In *Workshop on Word Sense Disambiguation: Recent Successes and Future Directions*. Philadelphia: ACL.
- Ide, Nancy and Jean Véronis, 1994. Multext (multilingual tools and corpora). In *Proceedings of the 15th International Conference on Computational Linguistics*. Kyoto.
- Krstev, Cvetana, Gordana Pavlović-Lažetić, Duško Vitas, and Ivan Obradović, 2004. Using Textual and Lexical Resources in Developing the Serbian Wordnet. *Romanian Journal of Information Science and Technology*, 7(1–2):147–162.
- Pavlović-Lažetić, Gordana, Duško Vitas, and Cvetana Krstev, 2004. Towards full lexical recognition. In *Text, Speech and Dialogue*, Lecture Notes in Artificial Intelligence. Berlin: Springer-Verlag. To appear.
- Silberztein, Max, 2000. *INTEX*. Masson.
- Sperberg-McQueen, C. M. and Lou Burnard (eds.), 1999. *Guidelines for Electronic Text Encoding and Interchange, Revised Reprint*. The TEI Consortium.
- Sperberg-McQueen, C. M. and Lou Burnard (eds.), 2002. *Guidelines for Electronic Text Encoding and Interchange, The XML Version of the TEI Guidelines*. The TEI Consortium.
- Tufiş, Dan, 1999. Tiered Tagging and Combined Language Model Classifiers. In Fredrik Jelinek and Elmar Noth (eds.), *Text, Speech and Dialogue*, number 1692 in Lecture Notes in Artificial Intelligence. Berlin: Springer-Verlag.
- Varadi, Tamas and Csaba Oravecz, 1999. Morphosyntactic Ambiguity and Tagset Design for Hungarian. In *Proceedings of the EACL-99 Workshop on Linguistically Interpreted Corpora (LINC-99)*. Bergen: ACL.
- Vitas, Duško and Cvetana Krstev, 2001. Intex and Slavonic Morphology. In *4es Journées INTEX*. Bordeaux. In print.
- Vitas, Duško and Cvetana Krstev, 2003. Composite tense recognition and tagging in serbian. In *Proceedings of the EACL 2003 Workshop on Morphological Processing of Slavic Languages*. Budapest.
- Vitas, Duško, Cvetana Krstev, Gordana Pavlović-Lažetić, and Ivan Obradović, 2003. An Overview of Resources and Basic Tools for Processing of Serbian Written Texts. In *Proc. of the Workshop on Balkan Language Resources, 1st Balkan Conference in Informatics*.