

# Literature and Aligned Texts

Duško Vitas<sup>1</sup> and Cvetana Krstev<sup>2</sup>

<sup>1</sup>Faculty of Mathematics, Studentski trg 16

<sup>2</sup>Faculty of Philology, Studentski trg 3

University of Belgrade

Belgrade 11000, Serbia

{vitas,cvetana}@matf.bg.ac.yu

## Abstract

In this paper we describe the structure of the aligned French-Serbian literary corpus and some results obtained in analyzing it. The alignment was performed on the subsentential level, which means that each text has been segmented into the <seg> elements in order to obtain in all cases the 1:1 alignment of the original and its translation. This means that in some cases the sentences were further segmented in two or more parts. The insertion of <seg> tags was initially done automatically, and obtained results were then manually validated. Some results obtained by analyzing this corpus are presented.

## 1 Introduction

One concise outline of the state of the art in the field of the aligned corpora was recently presented in (Tufis 06)<sup>1</sup>. Most of the presented corpora represent certain sublanguages, particularly of legal or administrative texts. Such choice is motivated either by the aims of the corpus development – terminology extraction, establishment of the translation equivalents, etc – or by the fact that the large amounts of this type of texts can be easily accessed in electronic form and the copyright issues for them can be easily settled.

The aligned corpora of literary texts are relatively rare, especially those where one of the languages is a Slavic language. One early attempt to develop the multilingual aligned text was the production of the aligned Plato's *Republic*, in the scope of TELRI project (Erjavec *et al.* 98). The produced multilingual aligned text included several bitexts where one of the languages in the pair was a Slavic language. This experience showed that the alignment task was by no means the easy one, especially if for a text with complex logical layout it was difficult to establish the original version. Namely, the differences of translations on some of the languages were significant, so the alignment process did not yield the convincing results. The multilingual aligned text was later produced in the scope of the same project for Orwell's

1984, which proved to be a more suitable text and was since used in many applications.

Probably the largest aligned corpus of literary Slavic texts consists primarily of the translations from English to different Slavic languages (Barentsen 05)<sup>2</sup>. After some modifications were done to the logical layout of the texts, in order to satisfy, for instance, the upper limit of the paragraph length, the texts were aligned on the level of paragraphs. The need of developing the aligned corpora of Slavic languages, especially South Slavic languages, was recently stressed in (Paskaleva & Pacovski 06) and illustrated by the aligned text of the novel *The Master and Margarita* by Mikhail Afanasievich Bulgakov (Russian/Bulgarian/Macedonian). An interesting discussion on problems of development of the corpus of literary texts aligned at the paragraph level were discussed in (Gelbukh *et al.* 06) and illustrated in the example of the English-Spanish corpus of novels. In our opinion, the context in the corpora aligned at the paragraph level is usually too wide for their effective and precise exploitation.

On the other hand, literary texts illustrate more precisely than any other texts translator's strategies, doubts and solutions. Also, aligned corpora of literary texts, to the contrary of corpora of some sublanguage, offer variety of possibilities for fruitful usage. Language learning, lexicography, and contrastive studies, are just a few examples of possible applications. More specific application would be to use one of the languages in the aligned text as a meta-language for the verification of the equivalent meanings in the other language. Such approach was proposed for the refinement of synsets (sets of synonymous literals) in Wordnets (Krstev *et al.* 04).

It should be noted that literary texts are, as a rule, precisely translated and offer the insight into the language and grammatical phenomena that

<sup>1</sup>see also: <http://www.cs.unt.edu/~rada/wa/>

<sup>2</sup>see also: <http://home.medewerker.uva.nl/a.a.barentsen/>

cannot be observed in other registers. Some of the examples are illustrated in this paper. The newspaper texts are often translated with greater freedom. One example is given in the fragment from the translation to Serbian of *Le Monde diplomatique* (see Appendix A). The third sentence in paragraph is omitted, and the adjective *tentaculaire* (engl. *tentacular*) from the French text is translated by the phrase *poput ogromne hobotnice sa dugim pipcima* (engl. *like the giant octopus with the long tentacles*). This metaphoric translation is further distributed through the whole paragraph (these occurrences are given in italic).

## 2 Corpus description

### 2.1 Sources

The corpus consists mostly of the classical works of French literature. The choice of texts was guided by their literary value, but also by their accessibility in electronic form. The largest part of French texts was obtained in electronic form from some web sites<sup>3</sup>, the other texts were scanned and corrected. The French originals span the time period from the end of the XVIII century to the present time, while all Serbian translation were done after 1926. They were all done by the most prominent translators from French to Serbian, and they are in accordance with the contemporary Serbian norm. The corpus also contains one text that is translated from Serbian to French (text no. 6 in Appendix B). Some Serbian texts were obtained by OCR, the others were re-typed, and one was obtained from the translator. The edited versions of Plato's *Republic* and Orwell's *1984*, according to (Erjavec *et al.* 98), were added to the corpus.

The format of the source French texts varied from the plain ASCII without any mark-up, to pdf, including some texts from *Gallica* that represented just the scanned images of the original texts. Most of the Serbian texts were in MS-Word format while others were ASCII texts with paragraph mark-up.

Corpus contains two Serbian translations of some texts. The inclusion of multiple translations was motivated by the possibility to investigate different variations in contemporary Serbian, in the first place on morphological level, by using French as a meta-language. The multiple translations exist for three texts:

- The corpus contains the translations of Voltaire's *Candide* from 1934 and 1964. The analysis of these translations showed that the more recent translation is the correction of the older translation from 1934.
- *The manuscript found in Saragosa* from Jan Potocki is translated in Serbian from the French abridged edition<sup>4</sup> and from the Polish integral edition<sup>5</sup>.
- Jules Verne's *The journey around the world in 80 days* is represented in corpus by two independent translations in Serbo-Croatian<sup>6</sup>.

These translations were independently aligned with French original, as well as between themselves. The text of Potocki's novel was aligned to the Polish electronic version of text<sup>7</sup>, that was corrected according to the paper edition, while several parts of the Verne's novels were aligned with a number of other languages as Greek, English, Spanish, and Slovene.

### 2.2 The size of corpus

The corpus contains the integral texts of Serbian translations and their French equivalents. In some cases, for instance for the *Anthology of French Fantastic Novels*, the Serbian translation of good quality was chosen first and then the corresponding original was sought for. The corpus contains the total of 14 texts written by 21 authors (see Appendix B for the full description). The precise data on the size of the corpus are given in Table 1. It can be noted that the French texts are longer than Serbian texts (approximately by 25%, if measured by number of tokens). This can be explained by the usage of articles in French, the omission of subject in Serbian, etc. On the other hand, the number of types is much bigger in Serbian due to the Serbian rich morphology.

### 2.3 Corpus preparation

Logical layout of texts was marked using the minimal set of tags: `<head>`, `<p>` and `<seg>`. Wherever possible, the pagination of the source text was retained as a comment. The bibliographic

<sup>4</sup>days from 1st to 14th, edited by Roger Caillois, Gallimard, 1958. see no. 12 in Appendix B.

<sup>5</sup>Jan Potocki: *Rękopis znaleziony w Saragossie*, Czytelnyk, Warszawa, 1965

<sup>6</sup>Belgrade 1962; Zagreb, 1961

<sup>7</sup>see <http://univ.gda.pl/~literat/sarag/index.htm>

<sup>3</sup>see <http://abu.cnam.fr/> and <http://gallica.bnf.fr/>

NO. OF NOVEL	TOKENS IN FRENCH	TYPES IN FRENCH	TOKENS IN SERBIAN	TYPES IN SERBIAN
1.	68843	10975	53666	13712
2.	27223	5191	22547	6021
3.	72253	9496	58661	11903
4.	93550	14662	78120	19385
5.	4710	1437	3743	1668
6.	39124	5132	31043	5774
7.	7158	1780	6083	2086
8.	33001	5457	29282	7617
9.	119489	10900	103127	17380
10.	960	472	769	509
11.	3654	1156	2957	1311
12.	65029	8178	51640	12024
13.	113419	11576	90026	16832
14.	154000	9266	108595	14184
$\Sigma$	802413	38247	640259	65162

Table 1: Corpus size (number of novel correspond to the numbers given in Appendix C)

description of the source text was given according to *TEI Guidelines*<sup>8</sup> in the TEI header element.

Tags `<head>` for the titles were inserted manually, while tags `<p>` for the paragraphs were inserted automatically or semi-automatically, depending on the format of the source text.

The recognition of sentence boundaries and insertion of `<seg>` tags was done by the finite-state transducer *Sentence.grf* that was implemented in the programming system *Intex*<sup>9</sup>. This transducer identifies the problems caused by the ambiguous usage of punctuation marks and with the high precision inserts a `<seg>`-tag at the beginning of the sentence. It enables, for instance, the distinction between the usage of the ASCII character 2E as a full stop from the other possible usages, like *Prof. X. Savary* in French or *7. 5. 2007.* in Serbian. This transducer is language dependent: for instance, a semi-colon in French is always at the end of the sentence, while this is not the case for Serbian.

All texts were aligned using the program *Xalign*<sup>10</sup>. The associated program *Concordancier* enables visualization and relinking of the concordances of the aligned texts and detection of the alignments that are not 1 : 1. Many-to-one correspondences can have various origins. The most common of them are:

- *Differences between an original and its translation.* For instance, the Serbian translation of Flaubert’s *Bouvard et Pécuchet* contains a draft of the end of the novel that is not

attached to the French text that served as a source. Thus, before the alignment process the Serbian text had 40 paragraphs more than the French original.

- *The omission of some parts of the text in translation.* For instance, the publisher obtained the text no. 9 (Appendix C) by scanning the previous edition. In this process one page was omitted by mistake, so the translation is shorter from the original by 14 sentences. The translators can make the similar mistake: in the same text seven sentences were not translated, in *Tours du monde en 80 jours* eight sentences, etc.
- *Differences in the paragraph marking.* The original and its translation often have different numbers of paragraphs. Although the paragraph represents the unit of the logical layout of one document, translators sometimes separate it into two or more paragraphs. For instance, if the original text contains sequence `...X said: </seg><seg>”...”` it can be transformed into `...X reče: </seg></p><p><seg> – ...` in the translation. Different document formats and various ways of their acquisitions can also lead to discrepancies between the number of paragraphs in the original text and its translation.
- *Differences in the sentence segmentation.* This type of many-to one correspondences is the most frequent one and it happens mostly because of the various orthographic issues that influence the segmentation process. For

<sup>8</sup>see <http://www.tei-c.org/P5/>

<sup>9</sup>see <http://intex.univ-fcomte.fr/>

<sup>10</sup>see <http://led.loria.fr/outils/ALIGN/align.html>

Paralelni prevod	
Francuski	Srpski
n1: En l'année 1872, la maison portant le numéro 7 de Saville-row, Burlington Gardens -- maison dans laquelle Sheridan mourut en 1814 --, était habitée par Phileas Fogg, esq., l'un des membres les plus singuliers et les plus remarquables du Reform-Club de Londres, bien qu'il semblât prendre à tâche de ne rien faire qui pût attirer l'attention.	n1: Godine 1872, u kući broj 7 u Ulici Sevil-rou Burlington Gardenz, u kojoj je 1816. godine umro Šeridan, stanovao je gospodin Fileas Fog, jedan od najčudnovatijih i najzapaženijih članova londonskog Reform-kluba, iako je izgledalo da se on trudi da ne učini ništa što bi moglo na njega privući pažnju.
n2: A l'un des plus grands orateurs qui honorent l'Angleterre, succédait donc ce Phileas Fogg, personnage énigmatique, dont on ne savait rien, sinon que c'était un fort galant homme et l'un des plus beaux gentlemen de la haute société anglaise.	n2: Posle jednog od najvećih govornika kojim se ponosi Engleska došao je, dakle, Fileas Fog, zagonetna ličnost o kojoj se nije ništa drugo znalo sem da je to čovek vrlo uglađen i jedan od najplemenitijih džentlmena visokog engleskog društva.
n3: On disait qu'il ressemblait à Byron -- par la tête, car il était irréprochable quant aux pieds --, mais un Byron à moustaches et à favoris, un Byron impassible, qui aurait vécu mille ans sans vieillir.	n3: Govorilo se da je licem sličan Bajronu - jer su mu noge bile bez zamerke -, ali jednom Bajronu sa brkovima i zaliscima, jednom hladnom Bajronu koji bi mogao živeti hiljadu godina a da ne ostari.
n4: Anglais, à coup sûr, Phileas Fogg n'était peut-être pas Londonner.	n4: Iako je na prvi pogled bio Englez, Fileas Fog verovatno nije bio Londonac.
n5: On ne l'avait jamais vu ni à la Bourse, ni à la Banque, ni dans aucun des comptoirs de la Cité.	n5: Niko ga nije nikada video na berzi, ili u banci, ili u jednoj od ustanova u centru grada.
n6: Ni les bassins ni les docks de Londres n'avaient jamais reçu un navire ayant pour armateur Phileas Fogg.	n6: Ni u pristaništa, ni u dokove Londona nije nikada prispeo nikakav brod čiji bi brodovlasnik bio Fileas Fog.
n7: Ce gentleman ne figurait dans aucun comité d'administration.	n7: Ovaj džentlmen nije se pojavljivao ni u jednom administrativnom nadležstvu.

Figure 1: The HTML visualization of the French-Serbian aligned text of *Tours du monde en 80 jours*

instance, in the following segment from the novel *Bouvard et Pécuchet*:

```
<p><seg> – "Mon Dieu, oui ! </seg>
<seg> On pourrait prendre le mien à mon
bureau!" </seg></p>
<p><seg> – Eh, bože, naravno, mogao bi
mi je ko uzeti u kancelariji! </seg></p>
```

the exclamation mark is used in the original, which is omitted in the translation, so there is no sentence ending at that point and the 2:1 correspondence is obtained. The opposite case can be found in the next segment:

```
<p><seg> – "Messieurs, je vous écoute !
quel est votre mal ? " </seg></p>
<p><seg> – Gospodo, slušam vas! </seg>
<seg> Na šta se žalite? </seg></p>
```

where the alignment is 1 : 2. Texts were manually edited in all such cases by insertion or omission of one pair of tags `</seg>...<seg>` in order to obtain the alignment 1 : 1 in all cases. At the same time, the validation of the whole alignment process was performed.

The results obtained by the described alternations in document's logical layout can be illustrated by some figures related to the novel *Bouvard et Pécuchet*: at the beginning, the initial French text had 3258 paragraphs (`<p>` tags) and 7163 sentences (`<seg>` tags), while the Serbian text had 3672 paragraphs and 6488 sentences. After the editing procedure both texts had 3108 paragraphs and 7075 sentences. Thus, it can be said that the alignment procedure aligns sub-sentential segments, that is, the `<seg>` tag is used to mark the equivalent segments that are whole sentences or parts of the sentence.

## 2.4 The corpus processing

The processing was performed in three steps. In the first phase all texts were independently processed using *Intex* as the processing tool which enabled their segmentation into sentences. Besides that, French and Serbian lexical resources integrated into the *Intex* system were applied to the

particular texts in order to obtain the insight into their lexical structure. As a byproduct, this procedure also enabled the correction of some spelling errors.

In the second phase, texts prepared in the first phase were separately aligned using the system Xalign. For each of the obtained bitexts the TMX format was produced using the system WS4LR (Krstev *et al.* 06), as well as HTML format that enables the visualization (see Figure 1).

In the third phase all texts were assembled and processed with the system IMS/CQP<sup>11</sup> that enabled the exploitation of the corpus as a whole.

### 3 Some examples

#### 3.1 The interjections

One of the peculiarities of the dialogs in the literary texts is the usage of interjections. We will consider only three French interjections: *ah!*, *eh!* *oh!* and their translations. Their role is interesting for two reasons. First, the interjections potentially mark the end of the sentences. Second, as all three interjections exist in Serbian as well, they could be regarded as cognates. The frequencies of their usage, however, does not support the validity of these suppositions: there are only 297 occurrences of these interjections in Serbian texts in comparison to 749 occurrences in French texts:

French	Serbian
<i>ah</i> 220	<i>ah</i> 173
<i>eh</i> 352	<i>eh</i> 73
<i>oh</i> 177	<i>oh</i> 51

From 352 occurrences of *eh*, 302 are in the expression *eh bien!*.

According to the French-Serbian (Serbo-Croatian) dictionary (Putanec 89), the suggested translations of these interjections are:

French	Serbian
<i>ah!</i>	<i>o(h)!</i> <i>jao!</i> <i>ah!</i>
<i>eh!</i>	<i>eh!</i> <i>ah!</i> <i>uh!</i> <i>aj!</i>
<i>eh bien!</i>	<i>no!</i> <i>dakle!</i>
<i>oh!</i>	<i>o!</i> <i>oh!</i> <i>ah!</i>

The concordances of the aligned corpus reveal the following equivalents:

French	Serbian
<i>ah!</i>	<i>ah!</i> <i>e!</i> <i>aha!</i> <i>ali!</i> <i>no!</i> <i>o!</i> <i>oh!</i> <i>a!</i>
<i>eh!</i>	<i>eh!</i> <i>he!</i> <i>o!</i>
<i>oh!</i>	<i>oh!</i> <i>eh!</i> <i>ah!</i> <i>ooo!</i>

<sup>11</sup><http://www.ims.uni-stuttgart.de/projekte/CorpusWorkbench/>

In particular, the interjection *et bien!* has been translated freely:

French	Serbian
<i>eh bien!</i>	<i>dakle!</i> <i>pa lepo!</i> <i>i onda?</i> <i>pa dobro!</i> <i>da, pa?</i> <i>pa onda?</i> <i>pa šta!</i> <i>eto vidiš!</i> <i>šta veliš!</i> <i>e pa lepo!...</i>

This example shows that the translation of interjections is often determined by translator's understanding of the context.

#### 3.2 The expression *sans doute*

Among the French expressions containing *doute*, dictionary (Robert 04) lists the expressions *sans doute* as the adverbial expression with description *selon toutes les apparences* (engl. *according to all the appearances*) and adds the comment that today this expression bears some doubt: *peut-être* (*maybe*), *probablement* (*probably*). On the contrary, the adverbial expression *sans nul* (*aucun*) *doute* has the meaning *certainement* (*without any doubt*). The Serbian equivalent *bez sumnje* according to the Serbian explanatory dictionary (Stevanović 76) has the meaning of French *sans aucun doute*. The bilingual French-Serbian dictionaries rarely indicate this subtle difference.

The analysis of the occurrences of the form *doute* in corpus gives the following results. In the French part of the corpus *doute* occurs 366 times, and almost all of them, except 58, are part of the phrases *sans (aucun) doute*. The French phrase *sans doute* occurs 270 times, while phrase *sans aucun doute* occurs 38 times. In the Serbian part of the corpus the noun *sumnja* occurs 143 times, and it is part of the phrase *bez sumnje* only 69 times.

The translator has to decide between two possibilities relying on his/her own intuition. The corpus search with the phrase *sans doute* retrieves the following equivalents in Serbian:

French	Serbian
<i>sans doute</i>	<i>verovatno, nesumnjivo, jamačno, van svake sumnje, sigurno, očigledno, bez sve sumnje, naravno, svakako,...</i>

Translation depends on a context, so the following equivalences can be spotted *aveugle sans doute* = *slepa za sve* (engl. *blind for everything*).

The phrase *sans aucun doute* is translated as *sigurno, bez sumnje, bez ikakve sumnje*.

If the equivalences are looked for in the only text that is translated from Serbian to French,

then the different set of possible equivalences is retrieved:

French	Serbian
<i>sans doute</i>	<i>valjda, verovatno, naročito, jesam, baš</i>
<i>sans aucun doute</i>	<i>bez sumnje</i>

Translation of the phrase *sans doute* shows that even the top-class translators, in this case to Serbian, can be misled to use the literal translation of the phrase that has the diachronic meaning. It is interesting to notice that in this case the dictionary can not be of much help since it is not clear when did the meaning of *sans doute* move from denoting certainty to probability.

### 3.3 Diminutives

In Serbian language diminutives and augmentatives can be produced in regular manner from most of the nouns (Vitas 04). Dictionaries, like the explanatory dictionary (Stevanović 76) record them only exceptionally, while in the bilingual lexicography there is no systematic way to describe this phenomenon that is on the border of the inflection and the derivation (Vitas & Krstev 04). Some of the diminutives that illustrate this phenomenon in the Serbian part of the corpus are *baroničica* besides *baronica* (engl. *baronette*), *Cigančica* but not *Ciganka* (engl. *Gipsy woman*), *divljačić* and also *divljak* (engl. *savage*), *kajganica* but not *kajgana* (engl. *scrambled eggs*) etc. In the Serbian part of corpus approximately 200 forms of diminutives occur. Corpus reveals two strategies for the usage of diminutives. The first one translates French sequence *petit N*, where *N* is a noun, with a form of diminutive, and the second one tries to fill the lexical gap for the concept that is not lexicalized in Serbian. The examples from the corpus for the first strategy are:

French	Serbian
<i>une petite veine</i>	<i>žilica</i>
<i>une petite bouteille</i>	<i>flašica</i>
<i>la petite bête</i>	<i>životinjica</i>
<i>les petites bricoles</i>	<i>stvarčice</i>
<i>le petit ruban</i>	<i>mašnica</i>
<i>les petites bouilles</i>	<i>čorbice</i>
<i>un petit banc</i>	<i>klupica</i>

Translators resort to diminutives also in cases when French lexeme contains some other modification of the basic lexeme, as in cases:

French	Serbian
<i>quelques gouttes</i>	<i>kapljice</i>
<i>faible portion</i>	<i>delić</i>
<i>la cordelette (&lt; le corde fine)</i>	<i>konopčić</i>
<i>le corbeille (&lt; le panier léger)</i>	<i>korpica</i>

The examples excerpted from the corpus that illustrate the second strategy are given in Table 2.

The listed examples show that the translators often resort to the usage of diminutives when they cannot find the more appropriate translation equivalent in Serbian.

### 3.4 The proper names

The proper names can be regarded as potential cognates that can be used for the correction of the results of the alignment (Krstev *et al.* 05) and (Vitas & Krstev 04). However, the translation of proper names from French to any language with rich morphological system generates three types of correspondence problems:

- Some problems arise from the possible use of two different alphabets (Latin and Cyrillic) and from the fact that in Serbian translation, the proper names are transcribed according to the phonetic orthography. For instance, *Bouvard* is transcribed as *Buvar*, *Passepartout* as *Paspartu*, etc., and that disables the direct usage of algorithms that are used by alignment software, as LCSR (Longest Common Subsequence Ratio) (Melamed 01) or Levenshtein distance.
- French personal names can occur in Serbian translation as a personal name and its possessive adjective in some of the inflective forms. The translation of a toponym is a toponym and its relational adjective. This phenomenon is illustrated by the following examples of the translations of proper names:

French	Serbian
<i>Bouvard</i>	<i>Buvar</i> realized as <b>N</b> : <i>Buvar, Buvara, Buvaru, Buvaru, Buvarom</i> and as <b>AdjPoss</b> : <i>Buvarov, Buvarova, Buvarove, Buvarovi, Buvarovih, Buvarovim, Buvarovo, Buvarovoj, Buvarovom, Buvarovu</i>
<i>Thérèse</i>	<i>Tereza</i> realized as <b>N</b> : <i>Tereza, Tereze, Terezi, Terezu, Terezo, Terezom</i> and as <b>AdjPoss</b> : <i>Terezine, Terezinih</i>
<i>Paris</i>	<i>Pariz</i> realized as <b>N</b> : <i>Pariz, Pariza, Parizom, Parizu</i> and as <b>AdjPoss</b> : <i>pariska, pariske, pariski, pariskim, pariskom</i>

	N	"PETIT" N	N+DEM
Sr	<i>brod</i>	<i>mali brod</i>	<i>brodić</i>
Fr	<i>le navire, le bateau, le paquebot</i>	<i>une goélette</i>	<i>une goélette</i>
Sr	<i>soba</i>	<i>mala soba</i>	<i>sobica</i>
Fr	<i>la chambre</i>	<i>une petite chambre</i>	<i>le cabinet, une petite pièce, l'arrière boutique</i>
Sr	<i>vrata</i>	<i>mala vrata</i>	<i>vratanca</i>
Fr	<i>la porte</i>	<i>une petite porte</i>	<i>les battants, une petite porte</i>
Sr	<i>čovak</i>	<i>mali čovak</i>	<i>čovečuljak</i>
Fr	<i>l'homme</i>	<i>un petit homme</i>	<i>le bonhomme, le petit homme</i>
Sr	<i>sto</i>	<i>mali sto</i>	<i>stočić</i>
Fr	<i>la table</i>	<i>la petite table</i>	<i>une console, un guéridon, une table de nuit</i>
Sr	<i>grad</i>	<i>mali grad</i>	<i>gradić</i>
Fr	<i>la ville</i>	<i>la petite ville</i>	<i>le petit bourg = un gros village!</i>

Table 2: In the column N+Dem, the examples of diminutives are given that are used to fill the lexical gap.

- In many cases, the frequencies of proper names are not equal in the original and its translation. The frequency data for the corresponding proper names in our corpus are:

French	Serbian
<i>Bouvard</i> 618	<i>Buvar.N</i> 617
<i>Phileas</i> 316	<i>Fileas.N</i> 314
<i>Fogg</i> 655	<i>Fog.N</i> 673
<i>Passepartout</i> 423	<i>Paspartu.N</i> 433
<i>Thérèse</i> 238	<i>Tereza.N</i> 238
<i>Paris</i> 106	<i>Pariz.N</i> 109

The difference in the number of occurrences do not necessarily mean that there actually exists the 1 – 1 correspondence in the number of cases that is equal to the smaller number of two frequencies. One example of the complex anaphora is given by the next segment:

```
<p><seg id="n2049"> Enfin, Bouvard et
Pécuchet s'adressèrent à Larsonneur.
</seg></p>
<p> ... </p>
<p> ... </p>
<p><seg id="n2053"> Par Gorju, ils s'en
procurèrent une douzaine, lui expédièrent
la moins grande – les autres enrichirent le
muséum. </seg></p>

<p><seg id="n2049"> Najzad se Buvar i
Pekiše obratiše Larsoneru. </seg></p>
<p> ... </p>
<p> ... </p>
<p><seg id="n2053"> Preko Goržija nabaviše
jedno tuce tih sekira, Larsoneru poslaše
najmanju, a ostalima obogatiše svoj
muzej.</seg></p>
```

The pronoun *lui* from the French segment n2053 refers to the personal name *Larsonneur* that was mention in the segment n2049. The translator, however, chose not to use the pronoun, believing perhaps that there would be a too big distance between the pronoun and its referent.

## 4 The conclusions

The first results in exploitation of the aligned French-Serbian corpus of literary texts indicate that in translations plenty of solutions can be found that are not recorded in bilingual French-Serbian dictionaries. Also, this corpus enables the analysis of translation strategies in resolving the lexical gaps or ambiguities in the original text, as well as discovering the sources of the inconsistencies in translation.

The further work in the aligned corpus exploitation might encompass the identification of the corresponding structures in the source and target language in the sense in which it has been suggested in (Blanco 01):

*Un système de TA doit donc être basé non sur des dictionnaires bilingues (ni, à plus forte raison, multilingues) mais sur [...] des descriptions lexicales de différentes langues effectuées d'après les mêmes principes.*

The existence of tools and resources in the same format for both languages, French and Serbian, makes the continuation of the experiments performed on this corpus feasible.

## 5 Appendix A: Non-literal translation – *Le Monde diplomatique* (May 2001)

### 5.1 French:

```
<head>La pieuvre publicitaire</head>
<p><seg id="n1"> Tentaculaire, étouffante, op-
pressive, la publicité ne cesse d'étendre ses do-
maines d'intervention. </seg> <seg id="n2">
Elle a récemment conquis de nouveaux terri-
toires, en particulier ceux de la galaxie Inter-
net. </seg> <seg id="n3"> Le chiffre d'affaires
publicitaire sur la Toile, en France l'an dernier,
avant la crise actuelle, a dépassé le milliard de
```

francs, soit plus que les recettes publicitaires des salles de cinéma. </seg> <seg id="n4"> Sous la forme discrète du parrainage, son champ d'intrusion ne connaît pratiquement plus de limites. </seg> <seg id="n5"> Par ce biais quasi clandestin, elle est parvenue à investir, ces dernières années, l'art, la culture, la science, l'éducation, et même la religion. </seg> </p>

## 5.2 Serbian:

```
<head>Čudovišna reklama</head>
<p><seg id="n1"> Poput ogromne hobotnice sa dugim pipcima, kojima će stegnuti i ugušiti sve što joj se nadje na putu, reklamna delatnost neumorno širi polje svog dejstva. </seg> <seg id="n2"> Nedavno je osvojila nova područja, medju kojima je posebno značajna Internet galaksija. </seg> <!-- missing segment --> <seg id="n3"> Njeni kraci dopiru svuda; najzad, sa ulogom sponzora koju je preuzela, reklama stiće neograničene mogućnosti. </seg> <seg id="n4"> Uvlači se poslednjih godina, neupadljivo i gotovo kradomice, u umetnost, kulturu, nauku, obrazovanje, čak i u religiju. </seg></p>
```

## 6 Appendix B: Corpus composition

1. *Anthology of French Fantastique* (edited by Zoran Mišić, Orfej 37, NOLIT, Beograd, 1964), including extracts from: J. Cazotte: *Le Diable amoureux*, D.A.F. de Sade: *Juliette*, Ch. Nodier: *Smarra*, G. de Nerval: *Aurelia*, H. de Balzac: *L'Élixir de longue vie* and *Séraphta*, P. Mérimée: *La Vénus d'Ille*, G. de Maupassant: *Horla, Qui sait?* and *La nuit*, Villiers de L'Isle-Adam: *L'Éve future*, *Véra* and *Le tueur de cygnes - Tribulat Bonhomet*, Comte de Lautramont: *Les Chants de Maldoror*, G. Apollinaire: *Le Roi-Lune*;
2. P. Louys: *La femme et le pantin*;
3. J. Verne: *Tours du monde en 80 jours*;
4. G. Flaubert: *Boward et Pécuchet*;
5. P. Besson: *Zodiaque amoureux*;
6. B. Blagojević: *Sve zveri što su sa tobom* (French translation: *L'arche de Boba*);
7. Voltaire: *Micromégas*;
8. Voltaire: *Candide ou l'optimisme*;
9. D.A.F. de Sade: *Justine*;
10. J.-P. Proudhon: *Idee générale de la Révolution au XIXe siècle* (extract);
11. F. Arrabal: *L'enterrement de la sardine*;
12. J. Potocki: *Manuscrit trouvé à Saragosse*;
13. G. Orwel: *1984*;
14. Platon: *La République: du régime politique*

## References

- (Barentsen 05) A. A. Barentsen. *Aspac – amsterdam slavic parallel aligned corpus*. Personal communication, Novi Sad, 2005.
- (Blanco 01) X. Blanco. Dictionnaires électroniques et traduction automatique espagnol-français. *Langages*, 143:143–166, 2001.
- (Erjavec et al. 98) T. Erjavec, A. Lawson, and L. Romary (Eds.). *East meets West A compendium of Multilingual Resources*. TELRI Association e.V., IdS, Mannheim, 1998.
- (Gelbukh et al. 06) A. Gelbukh, G. Sidorov, and J.A. Vera-Felix. A bilingual corpus of novels aligned at paragraph level. In T. Salakoski and F. Ginter, S. Pyysalo, and T. Pahikkala, editors, *Advances in Natural Language Processing, FINTAL 2006*, LNAI, 4139, pages 16–23. Springer Verlag, 2006.
- (Krstev et al. 04) C. Krstev, D. Vitas, G. Pavlović-Lažetić, and I. Obradović. Using textual and lexical resources in developing serbian wordnet. *Romanian Journal of Information Science and Technology*, 7:147–161, 2004.
- (Krstev et al. 05) C. Krstev, D. Vitas, D. Maurel, and M. Tran. Multilingual ontology of proper names. In Z. Vetulani, editor, *Proceedings of 2nd Language and Technology Conference*, pages 116–119. Wydawnictwo Poznańskie Sp. z o.o., Poznań, 2005.
- (Krstev et al. 06) C. Krstev, R. Stanković, D. Vitas, and I. Obradović. Ws4lr - a workstation for lexical resources. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation, LREC'06*, pages 1692–1697. Genoa, Italy, 2006.
- (Melamed 01) D. Melamed. *Empirical Methods for Exploiting Parallel Texts*. The MIT Press, Cambridge, MA, London, 2001.
- (Paskaleva & Pacovski 06) E. Paskaleva and I. Pacovski. Aligning the translation – a possible strategy for creation of aligned corpora (for south-slavic languages). In S. Koeva and M. Dimitrova-Vulchanova, editors, *The Fifth International Conference "Formal Approaches to South Slavic and Balkan Languages", FASSBL V*, pages 113–117. Institute of Bulgarian Language, Bulgarian Academy of Sciences, Sofia, 2006.
- (Putanec 89) V. Putanec. *Dictionnaire français-croate ou serbe*. Školska knjiga, Zagreb, 1989.
- (Robert 04) P. Robert. *Dictionnaire Le Petit Robert*. Le Robert, Paris, 2004.
- (Stevanović 76) M. Stevanović. *Rečnik srpskohrvatskoga književnog jezika*. Matica srpska : Matica hrvatska, Novi Sad, Zagreb, 1967–1976.
- (Tufis 06) D. Tufis. Cross-lingual knowledge induction from parallel corpora. In S. Koeva and M. Dimitrova-Vulchanova, editors, *The Fifth International Conference "Formal Approaches to South Slavic and Balkan Languages", FASSBL V*, pages 15–18. Institute of Bulgarian Language, Bulgarian Academy of Sciences, Sofia, 2006.
- (Vitas & Krstev 04) D. Vitas and C. Krstev. Intex and aligned texts. In C. Muller, J. Royaut, and M. Silberstein, editors, *Intex pour la linguistique et le traitement automatique des langues*, pages 255–270. Presses Universitaires de Franche Comté, Besançon, 2004.
- (Vitas 04) D. Vitas. Morphologie dérivationnelle et mots simples: Le cas du serbo-croate. In C. Leclere, E. Laporte, M. Piot, and M. Silberstein, editors, *Lexique, Syntaxe et Lexique-Grammaire / Syntax, Lexis and Lexicon-Grammar (Papers in honour of Maurice Gross)*, volume 24 of *Linguisticae Investigationes Supplementa*, pages 629–640. John Benjamin, Amsterdam/Philadelphia, 2004.