

Automatic Translation Lexicon Extraction from Czech-English Parallel Texts

Jan Cuřín, Martin Čmejrek

Abstract

We present experimental results of an automatic extraction of a Czech-English translation dictionary by means of procedures based on a probabilistic approach. We used two different bilingual corpora (computer oriented of 119,886 and fiction of 58,137 sentence pairs). For the automatic sentence alignment a statistical method was used based on sentence lengths, for the dictionary extraction a regular grammar based noun group extractor and the probabilistic model of (Brown *et al.*) were combined. The size of the resulting dictionaries is around 6,000 entries. After the significance filtering, weighted precision is 86.4% for the computer oriented Czech-English dictionary and 70.7% for the fiction part.

1 Introduction

The primary motivation for our research was to create translation lexicon of terminology of a particular discipline. Many disciplines lack relevant dictionaries or the dictionaries are obsolete because of the quick development of the discipline. The idea was that the basic part of the translation lexicon would be generated from the parallel corpus of the hitherto translated texts automatically and afterwards it could be manually edited.

The works in the field of automatic sentence alignment (Gale and Church, 1993) and automatic extraction of translation dictionary (Brown *et al.*, 1993; Wu and Xia, 1994) have exploited very large corpora. The former two used Canadian Hansards English-French Corpus, the latter one used the HKUST English-Chinese Corpus. These corpora include speeches delivered in bilingual parliaments in Canada and Hongkong. They are highly equivalent and satisfy rather strict constraints put on parallel corpora. The translations are mostly literal and tight. The situation in our country is different, we lack such a good source of large bilingual data.

Therefore, we decided to use a smaller corpus of texts taken from a particular discipline - a computer oriented corpus. This corpus consists of operating system messages from IBM AIX and of operating system guides for IBM AS/400 and VARP 4. These data are products of localization and translation of software from English into Czech. The translations are very literal and precise. In most cases sentences are translated sentence by sentence, it means that there is a one-to-one correspondence between an English sentence and a Czech sentence. But on the other hand, it is a typical feature of this kind of texts that majority of operating system messages and a big part of sentences from guides are nominal sentence, i.e. they do not have a verb.

We have also got an access to data from the Reader's Digest Výběr magazine. Every issue of this global magazine contains 30-60 % of articles that have been translated from English to the local language. We had to search in the English version to find the corresponding articles that are in the Czech version. There was also a lot of manual work to do because of the various word processor formatting. The translations in Reader's Digest are mostly very free. They include many constructions with direct speech. Some articles had to be excluded, for example such articles as "what to eat if we want to lose

weight” in which the American food was substituted by Central European meals. And such substitution concerned all culture-specific facts in the process of localization. We also excluded jokes as the majority of them was translated using completely different words and strategies.

The reason why we decided to carry out the experiments also with Reader’s Digest was to compare the results of methods applied to computer oriented texts on the one hand with those of journalistic texts and fiction on the other.

In section 2 we describe implementation of the statistical method for automatic paragraph and sentence alignment (Gale and Church, 1993). The important differences in the distribution of alignment categories between Canadian Hansards and Czech-English Fiction Corpora are documented.

Majority of the terms in our training data occur in the form of a noun group. The simple regular grammar based tool described in section 3 marks noun group boundaries. This is to our knowledge the first use of this method.

Section 4 describes probabilistic models 1 and 2 of (Brown *et al.*, 1993) and their implementation of the training procedure.

Output of the training procedure is filtered to produce a smaller, more useful dictionary. Section 5 is dedicated to this problem. The final evaluation and an example of the resulting dictionaries extracted from both corpora are presented.

2 Statistical Alignment of Paragraphs and Sentences in Czech-English Bilingual Corpora

For the subsequent training procedures we need to identify matching paragraphs and sentences between both languages automatically. There are two main approaches to this problem: a lexical and a statistical one, and many works use one of them or combine both (Church, Gale, 1993; Wu, Xia, 1994). Lexical approaches are based on on-line bilingual dictionaries while the statistical ones use simple probabilistic models usually based on lengths of aligned sentences. There is no on-line Czech-English dictionary available, so we have implemented the statistical model from (Gale and Church, 1993).

The problem can be formalized as follows:

Let us have Czech and English texts (typically paragraphs) T_c and T_e . The alignment is a set of pairs of parts of texts (typically 0, 1, 2 . . . m sentences) $\{(L_{c,1} \rightleftharpoons L_{e,1}), \dots, (L_{c,n} \rightleftharpoons L_{e,n})\}$ such that $L_{c,1}, \dots, L_{c,n} = T_c$ and $L_{e,1}, \dots, L_{e,n} = T_e$. We are looking for the best alignment \mathcal{A} that maximizes the probability over all possible alignments on T_c and T_e :

$$\arg \max_{\mathcal{A}} \Pr(\mathcal{A}|T_c, T_e). \tag{1}$$

Now we have to make several approximations to obtain an effectively computable model: First, the set of possible types of matching parts of texts is limited to six categories: 1-1, 0-1, 1-0, 1-2, 2-1, 2-2. We also assume that the probabilities of individual aligned parts of texts $\Pr(L_{c,i} \rightleftharpoons L_{e,i})$ are independent. And finally we assume that the probabilities of individual alignments depend only on the function δ of lengths of aligned parts l_c, l_e :

$$\arg \max_{\mathcal{A}} \Pr(\mathcal{A}|T_c, T_e) \approx \arg \max_{\mathcal{A}} \prod_{(L_c \rightleftharpoons L_e) \in \mathcal{A}} \Pr(L_c \rightleftharpoons L_e | \delta(l_c, l_e)). \tag{2}$$

Type of Alignment	Computer # sent.	Oriented $\Pr(L_c \bowtie L_e)$	Fiction # sent.	$\Pr(L_c \bowtie L_e)$	Canadian # sent.	Hansards $\Pr(L_c \bowtie L_e)$
1-1	109	0.90	64	0.69	1167	0.89
1-0 & 0-1	3	0.03	3	0.03	13	0.01
1-2 & 2-1	7	0.06	24	0.26	117	0.09
2-2	1	0.01	2	0.02	15	0.01
total	120	1.00	93	1.00	1312	1.00

Table 1: Sentence alignment type distribution on a hand-annotated sample.

Function δ is defined as follows:

$$\delta(l_c, l_e) = \frac{l_e - l_c e}{\sqrt{l_c \sigma^2}}, \quad (3)$$

where $e = E(r) = E(\frac{l_e}{l_c})$ is the mean number of English characters generated by each Czech character, and $\sigma^2 = D^2(\frac{l_e}{l_c})$ is the variance. After applying Bayes' Rule and with the assumption that δ is normally distributed:

$$\begin{aligned} & \arg \max_A \prod_{(L_c \rightleftharpoons L_e) \in A} \Pr(L_c \rightleftharpoons L_e | \delta(l_c, l_e)) \approx \\ & \approx \arg \min_A \sum_{(l_c \rightleftharpoons l_e) \in A} -\log \frac{\Pr(\delta(l_c, l_e) | l_c \rightleftharpoons l_e) \Pr(L_c \rightleftharpoons L_e)}{\Pr(\delta(l_c, l_e))} \approx \\ & \approx \arg \min_A \sum_{(L_c \rightleftharpoons L_e) \in A} -\log \left(\frac{2(1 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{|\sigma|} e^{-\frac{z^2}{2}} dz) \Pr(L_c \rightleftharpoons L_e)}{\Pr(\delta(l_c, l_e))} \right). \end{aligned} \quad (4)$$

Parameters e , σ^2 and $\Pr(L_c \rightleftharpoons L_e)$ are estimated from a sample of hand-aligned sentences. $\Pr(\delta(l_c, l_e))$ is constant and as such doesn't influence the result of minimization.

Table 1 compares parameters of three different corpora. Although Canadian Hansards and the Czech-English computer oriented corpus have very similar distribution of categories of alignment, the distribution differs substantially on the fiction corpus.

The best alignment is now easy to find if we use the dynamic programming procedure. In the first step, the algorithm aligns paragraphs in matching articles, in the second step, it aligns sentences in matching paragraphs.

Table 2 briefly summarizes results of the automatic paragraph and sentence alignment. The accuracy was 96 correctly aligned pairs on computer oriented corpora and 85 pairs on fiction corpora in a randomly selected sample of 100 pairs of sentences.

3 Groups Identification

We aim at a terminological dictionary, i.e. a dictionary containing also translations consisting of more than one word. For example, single word *typewriter* (in English) is translated by a noun group consisting in two words *psací stroj* (lit. 'writing machine') in Czech. On the other hand, the English group *construction worker* corresponds to a single Czech word *stavbař*. We model a tool which is capable of handling such cases.

	Computer Oriented	Fiction
# words (English)	1245780	959583
# words (Czech)	1089813	860757
# paragraphs (English)	88790	19567
# paragraphs (Czech)	88790	24874
# sentences (English)	120743	70872
# sentences (Czech)	121295	67856
# aligned sentences	119886	67436
types of alignment:		
1-1	117450	37039
0-1 (En/Cz)	73	5311
1-0 (En/Cz)	36	4454
1-2 (En/Cz)	1397	9501
2-1 (En/Cz)	882	7342
2-2	48	1089

Table 2: Results of automatic paragraph and sentence alignment.

The idea is to concatenate words of potential groups into one string, i.e. to consider these constructions as single "words". Identification of groups is based on a simple regular grammar. Grammar rules can be modified by the user. The grammar used in our case identifies noun groups (word sequences which consist from nouns, adjectives and some auxiliary words). Only continuous sequences of words are considered to be noun groups.

Czech is an inflective language with a lot of word forms, and in English contrastively, one word form corresponds to more POS categories. That is why we decided to proceed as follows:

- Czech groups and words are converted into their basic forms (the basic form means the first case of singular or plural for nouns, adjectives and pronouns, and the infinitive for verbs),
- English nouns and adjectives are distinguished from other POS categories,
- definite and indefinite articles are removed from English groups (there is no equivalent for the category of an article in Czech).

For an example, see Figure 1.

Tagging of Czech and English texts and conversion of Czech word forms were done by the BH tagging tools (Hajič, Hladká, 1998).

Marking of potential groups in a sentence is done separately for each language. Noun group identification algorithm works in two passes through the sentence. All possible groups in the sentence are identified in the first pass. During the second pass the algorithm searches for such combination of groups that:

- do not overlap in one combination,
- cover the maximal number of words in the sentence,
- the number of groups in the combination is minimal.

basic form	original forms
<i>integrováný systém souborů</i>	← <i>integrováný systém souborů</i>
	← <i>integrovaného systému souborů</i>
	← <i>integrovanému systému souborů</i>
	← <i>integrováný systéme souborů</i>
	← <i>integrovaném systému souborů</i>
	← <i>integrovaným systémem souborů</i>
<i>integrated file system</i>	← <i>the integrated file system</i>
	← <i>integrated file system</i>
	← <i>an integrated file system</i>

Figure 1: Conversion of groups into their basic form.

If there is still more than one such combination, one of them is chosen randomly. Parallel sentences with marked groups are shown in Figure 2. Concatenated noun groups or nouns (one word groups) are delimited by symbols **&** and **#**.

Thanks to the fact that the learning of the dictionary is based on probabilistic methods, we have a discretion in group identification. Even if some groups are marked incorrectly, they are often eliminated by the probabilistic algorithm which handles a big amount of mostly good data.

Once the data are prepared, the translation dictionary training procedure can start.

4 Translation Dictionary Training

We implemented models 1 and 2 described in (Brown *et al.*, 1993) of sentence translation probability, and used iterative EM algorithm for maximizing the likelihood of generating the Czech translation from the English text.

Basic definitions:

Let $\mathbf{c} \equiv \mathbf{c}_1, \dots, \mathbf{c}_{l_c}$ and $\mathbf{e} \equiv \mathbf{e}_1, \dots, \mathbf{e}_{l_e}$ be Czech and English sentences, and the word alignment be $\mathbf{a} \equiv \mathbf{a}_1, \dots, \mathbf{a}_{l_c}$, such that $0 \leq a_i \leq l_e$, represents the information that the Czech word c_i is a translation of English word e_{a_i} . Symbols c_0 and e_0 stand for an empty word. The training corpus S is a set of pairs (\mathbf{c}, \mathbf{e}) .

Model 1 in (Brown *et al.*, 1993) is based on word-by-word translation probability $t(c|e)$ and approximates the probability of translating the English sentence \mathbf{e} into the Czech sentence \mathbf{c} following a word alignment \mathbf{a} . This model also assumes, that all the possible word alignments are equally probable and also that all the possible lengths of Czech sentences have the same probability ϵ :

$$\Pr(\mathbf{c}, \mathbf{a}|\mathbf{e}) = \frac{\epsilon}{(l_e + 1)^{l_c}} \prod_{i=1}^{l_c} t(c_i|e_{a_i}). \quad (5)$$

The re-estimation formulas for EM algorithm are:

$$\bar{c}(c|e; \mathbf{c}, \mathbf{e}) = \frac{t(c|e)}{t(c|e_0) + \dots + t(c|e_{l_e})} \sum_{i=1}^{l_c} \sum_{j=0}^{l_e} \delta(c, c_j) \delta(e, e_i), \quad (6)$$

English: The device driver indicates a hardware failure of equipment.

`&_device_driver_#` indicates `&_hardware_failure_of_equipment_#` .

Czech: Ovladač zařízení zjistil technickou závadu přístroje.

`&_ovladač_zarizení_#` zjistit `&_technická_závada_přístroje_#` .

English: Just then, they saw cowboys coach Eddie Sutton walk toward the court with a man pushing a kid in a wheelchair.

just then they saw `&_cowboys_coach_#` eddie sutton walk toward `&_court_#` with `&_man_#` pushing `&_kid_#` in `&_wheelchair_#` .

Czech: V tu chvíli zahlédli, jak na hřiště přichází jejich trenér Eddi Sutton s mužem, který na invalidním vozíku vezl malého chlapce.

ten `&_chvíle_#` zahlédnout jak na `&_hřiště_#` přicházet jeho `&_trenér_#` eddi sutton s `&_muž_#` který na `&_invalidní_vozík_#` vzt `&_malý_chlapec_#` .

Figure 2: Sentences with marked groups.

$$\lambda_e = \sum_c \sum_{(\mathbf{c}, \mathbf{e}) \in S} \bar{c}(c|e; \mathbf{c}, \mathbf{e}), \quad (7)$$

$$t(c|e) = \sum_{(\mathbf{c}, \mathbf{e}) \in S} \frac{\bar{c}(c|e; \mathbf{c}, \mathbf{e})}{\lambda_e}. \quad (8)$$

Model 2 in (Brown *et al.*, 1993) is an extension of model 1 by using word alignment probabilities $a(a_i|i, l_{\mathbf{c}}, l_{\mathbf{e}})$:

$$\Pr(\mathbf{c}, \mathbf{a}|\mathbf{e}) = \epsilon \prod_{i=1}^{l_{\mathbf{c}}} t(c_i|e_{a_i})a(a_i|i, l_{\mathbf{c}}, l_{\mathbf{e}}). \quad (9)$$

The re-estimation formulas for the EM algorithm are:

$$\bar{c}(c|e; \mathbf{c}, \mathbf{e}) = \sum_{i=1}^{l_{\mathbf{c}}} \sum_{j=0}^{l_{\mathbf{e}}} \delta(c_i, c)\delta(e_j, e) \frac{t(c|e)a(j|i; l_{\mathbf{c}}, l_{\mathbf{e}})}{\sum_{j'=0}^{l_{\mathbf{e}}} t(c_i|e_{j'})a(j'|i; l_{\mathbf{c}}, l_{\mathbf{e}})}, \quad (10)$$

$$\bar{c}(i|j; l_{\mathbf{c}}, l_{\mathbf{e}}; \mathbf{c}, \mathbf{e}) = \sum_{i=1}^{l_{\mathbf{c}}} \sum_{j=0}^{l_{\mathbf{e}}} \delta(c_i, c)\delta(e_j, e) \frac{t(c|e)a(j|i; l_{\mathbf{c}}, l_{\mathbf{e}})}{\sum_{j'=0}^{l_{\mathbf{e}}} t(c_i|e_{j'})a(j'|i; l_{\mathbf{c}}, l_{\mathbf{e}})}, \quad (11)$$

$$\lambda_e = \sum_c \sum_{(\mathbf{c}, \mathbf{e}) \in S} \bar{c}(c|e; \mathbf{c}, \mathbf{e}), \quad (12)$$

$$\lambda_{i, l_{\mathbf{c}}, l_{\mathbf{e}}} = \sum_{0 < j < l_{\mathbf{e}}} \sum_{(\mathbf{c}, \mathbf{e}) \in S} \bar{c}(i|j; l_{\mathbf{c}}, l_{\mathbf{e}}), \quad (13)$$

$$t(c|e) = \sum_{(\mathbf{c}, \mathbf{e}) \in S} \frac{\bar{c}(c|e; \mathbf{c}, \mathbf{e})}{\lambda_e}, \quad (14)$$

$$t(i|j; l_{\mathbf{c}}, l_{\mathbf{e}}) = \sum_{(\mathbf{c}, \mathbf{e}) \in S} \frac{\bar{c}(i|j; l_{\mathbf{c}}, l_{\mathbf{e}})}{\lambda_{i, l_{\mathbf{c}}, l_{\mathbf{e}}}}. \quad (15)$$

The EM algorithm works as follows:

1. Set the consistent initial values for $t(\cdot) > 0$ and $a(\cdot) > 0$.
2. Compute all $c(\cdot)$.
3. Compute the Lagrange multiplier λ_e for each English word e and λ_{i,l_c,l_e} for each position of word word alignment $\langle i, l_c, l_e \rangle$.
4. Re-estimate $t(\cdot)$ using equation 8 for model 1 or equations 14 and 15 for model 2.
5. Repeat steps 2–4 until $t(\cdot)$ converge.

Subcorpus of sentences of type 1–1 was selected from the computer oriented corpus. Subcorpus of 1–1, 1–2, 2–1 and 2–2 sentences was selected from the fiction corpus.

5 Significance Filtering and the Evaluation of Results

The training procedure described in the previous section results in a probabilistic dictionary which assigns translation probability to every pair of Czech and English words, which have ever been seen together in corresponding sentences. It is necessary to "clean up" the probabilistic dictionary by filtering out most of the translations to produce a useful dictionary. An obvious solution to reduce translations is to set a threshold on probabilities. Absolute thresholds work poorly, we use them only for rough pruning of translations with negligible probabilities.

The principle for significant filtering is to find a combination of just a few filtering criteria that affects the quality of the representative sample of the dictionary in the best way. This combination is used to filter the whole dictionary.

At the beginning we set the dictionary quality indicators and manually mark the quality of translations in the representative sample of the dictionary. We use two obvious quality indicators Precision and Recall, and a third one, Share, defined as follows:

Let \mathcal{T} be a set of all translations in the input dictionary, \mathcal{G} be a set of good translations for corresponding entries (i.e. these translations were marked as good by hand in the representative part of the dictionary), and \mathcal{S} be a set of translations which were marked as good by the combination of filtering criteria (i.e. these translation were marked as good by the automatic filtering). We denote the translation probability of a particular translation (x) by $\text{Pr}(x)$. Let \mathcal{S}^* and \mathcal{G}^* be similar sets including only very good translations. Good translations in contrast to very good ones are acceptable only in some context.

$$\text{Precision}(\mathcal{T}) = \frac{\text{card}(\mathcal{S} \cap \mathcal{G})}{\text{card}(\mathcal{S})}, \quad (16)$$

$$\text{Share}(\mathcal{T}) = \frac{\sum_{x \in \mathcal{S} \cap \mathcal{G}} \text{Pr}(x)}{\sum_{x \in \mathcal{S}} \text{Pr}(x)}, \quad (17)$$

$$\text{Recall}^*(\mathcal{T}) = \frac{\text{card}(\mathcal{S}^* \cap \mathcal{G}^*)}{\text{card}(\mathcal{G}^*)}. \quad (18)$$

The quality indicator Precision grows with the elimination of bad translations. Share is a weighted Precision, which takes into account probabilities assigned to translations by the training procedure. Recall indicates a success in recognizing manually marked good translations by the automatic filtering.

The representative sample of the dictionary was about 4% of all entries.

Filtering criteria used are defined in Table 3.

critereon	description
Frst (n)	First n translations selected.
Thd (p)	Only translations accounting for the top of the threshold p are retained.
MC (n)	Works only for entries with low occurrence. Translations having count higher than n are excluded, if they have not been selected as groups. Translation probabilities for each entry are recomputed.
MPr (p)	Translations with the translation probability lower than p are excluded. Translation probabilities for each entry are recomputed.
NonT (p)	Works only for entries selected as groups. Translations with translation probabilities lower than p are excluded, if they have not been selected as groups. Translation probabilities for each entry are recomputed.

Table 3: Definiton of Filtering criteria.

<i>Combination of Criteria</i>	Computer Oriented R*/P/S (in %)	Fiction R*/P/S (in %)
Thd(0.85) \sim input dictionary	100.0/38.0/63.6	100.0/10.5/26.0
Frst(1)	56.6/88.1/88.1	57.3/61.7/61.7
Thd(0.7)	94.1/49.8/69.8	96.7/13.6/30.4
Thd(0.7) \rightarrow MPr(0.08)	82.4/60.6/73.1	73.4/47.3/55.9
Thd(0.7) \rightarrow MPr(0.05) \rightarrow MPr(0.09)	81.3/62.1/73.8	81.3/38.7/49.7
Thd(0.7) \rightarrow MPr(0.05) \rightarrow MPr(0.09) \rightarrow NonT(0.3)	81.3/68.5/78.4	80.9/41.7/54.4
MC(1800) \rightarrow Thd(0.7) \rightarrow MPr(0.05) \rightarrow MPr(0.09) \rightarrow NonT(0.3)	84.2/72.6/83.7	81.7/51.5/63.6
Thd(0.7) \rightarrow MC(1800) \rightarrow MPr(0.05) \rightarrow MPr(0.09) \rightarrow NonT(0.3)	83.8/74.8/84.9	81.3/51.4/63.9
Thd(0.7) \rightarrow MC(1200) \rightarrow MPr(0.05) \rightarrow MPr(0.09) \rightarrow NonT(0.3)	82.4/75.2/85.1	81.7/53.0/65.9

Table 4: Criteria combination and their influence on dictionary quality indicators.

	Recall*	Precision	Share
	Cz-En / En-Cz	Cz-En / En-Cz	Cz-En / En-Cz
<i>Computer Oriented</i>	86.8 / 83.8	75.1 / 74.8	86.4 / 84.2
<i>Fiction</i>	83.8 / 81.7	54.6 / 53.0	70.7 / 65.9
<i>Computer Or. (noun groups only)</i>	93.4 / 94.6	82.5 / 78.4	91.0 / 86.9
<i>Fiction (noun groups only)</i>	83.2 / 74.7	62.6 / 69.4	76.1 / 76.0

Table 5: Quality indicators in the output dictionaries.

Applying miscellaneous combinations of filtering criteria (changing the order of filtering criteria or criterion thresholds) we observe progress in dictionary quality indicators. An example can be seen in Table 4, where we show dictionary quality rates for English-Czech dictionaries (Recall* / Precision / Share). On the first line there are the rates for input dictionary (that is 100% for Recall*). The combination of filtering criteria with balanced values of dictionary quality indicators is chosen as optimal (the last row in Table 4). The whole dictionary is processed by this optimal combination of filtering criteria. Each output dictionary (Cz-En and En-Cz for computer and fiction) contains about 6.000 entries (see quality indicators in Table 5).

Figures 3 and 4 are examples of the filtered dictionaries. Entries and translations marked by * were recognised as groups by our noun group extracting tool. Numbers in square brackets are frequency counts in the corpus. Numbers in round brackets are translation probabilities (normalised for each entry). Good translations are underlined.

In Figure 3 there is a sample of English-Czech computer oriented lexicon. For instance entries **map** and **map*** or **mark** and **mark*** distinguish verbs and nouns translations. An example of a common error is the translation of the entry **manual IPL***, where the noun group was not recognised in the corresponding Czech sentence.

The frequency of good translations in the sample of Czech-English fiction translation lexicon in Figure 4 is lower than in the previous sample. The reason is the richness of tokens in the fiction part of corpora: for the same amount of data there were three times more tokens than in the computer oriented one. The same error of the non-recognised noun group appears for entry **rentgen*** and its translation *X-rays*. Just for fun realise that the entry **režisér*** (director) is more probable translated as *Spielberg* than to the exact translation *director*.

6 Conclusion

The reported experiments are, to our knowledge, the first demonstration of methods mentioned above for Czech and English parallel corpora. The results of automatic paragraph and sentence alignment on the computer oriented corpora reach a similar quality (96%) as those achieved on Canadian Hansards. Results on fiction corpora are worse (85%) because of the lower quality and non-literality of translations. The results of the dictionary extraction, for the computer oriented corpora are of unexpectedly high share (weighted precision) rates about 85% and for the terminology dictionary (that contains only noun groups) they are even better: 87%–91%. Soon, the results of this work will be used in practice for translation purposes.

manage [177] <u>spravovat</u> (0.47), <u>řídít</u> (0.31), <u>správa</u> * (0.22)	mapped [22] <u>mapovat</u> (1.00)
managed [21] <u>řídít</u> (0.36), <u>spravovat</u> (0.27), <u>program</u> (0.18), <u>server/400</u> (0.18)	mapping * [45] <u>mapování</u> * (0.45), <u>macintosh</u> (0.30), <u>přirazení</u> * (0.25)
management * [37] <u>management</u> * (0.78), <u>řízení</u> * (0.22)	maps [19] <u>mapy</u> * (0.56), <u>instalační</u> (0.22), <u>jeho</u> (0.22)
manager's maintenance operating * [10] <u>operating</u> (0.77), <u>podrobnější informace</u> * (0.23)	maps * [10] <u>mapy</u> * (0.62), <u>aplikace</u> * (0.38)
manager * [76] <u>manager</u> * (1.00)	margins * [13] <u>okraje</u> * (0.85), <u>řádek</u> * (0.15)
manager software operating * [13] <u>operating</u> (0.34), <u>nalézt</u> (0.32), <u>SC19</u> (0.18), <u>program</u> * (0.16)	mark [19] <u>označit</u> (1.00)
manages [14] <u>řídít</u> (1.00)	mark * [18] <u>označit</u> (0.54), <u>značka</u> * (0.46)
managing [87] <u>řízení</u> * (0.36), <u>správa</u> * (0.27), <u>spravující stroje</u> * (0.22), <u>spravující stroj</u> * (0.16)	marked [43] <u>označit</u> (0.83), <u>označený</u> (0.17)
managing system * [13] <u>řídící systém</u> * (1.00)	marketing representative * [62] <u>obchodní zástupce</u> * (1.00)
manual * [105] <u>manual</u> * (0.44), <u>manuál</u> * (0.36), <u>příručka</u> * (0.21)	marks [13] <u>uvést</u> (0.40), <u>klíčové slovo</u> * (0.40), <u>uvozovky</u> * (0.20)
manually [130] <u>ručně</u> (0.79), <u>manuálně</u> (0.21)	mask * [35] <u>maska</u> * (0.59), <u>maska podsítě</u> * (0.41)
manuals * [11] <u>příručky</u> * (0.57), <u>knihovna</u> * (0.21), <u>vyhledávání informací</u> * (0.21)	master * [13] <u>master</u> * (1.00)
manual installation * [13] <u>ruční instalace</u> * (1.00)	master installation list * [50] <u>hlavní instalační formulář</u> * (1.00)
manual installation process * [22] <u>proces ruční instalace</u> * (0.44), <u>ruční instalace</u> * (0.41), <u>proces</u> * (0.15)	match [177] <u>odpovídat</u> (0.87), <u>souhlasit</u> (0.13)
manual install display * [10] <u>obrazovka manual install</u> * (0.60), <u>objevit</u> (0.40)	match * [46] <u>odpovídat</u> (0.41), <u>odpovídající protějšek</u> * (0.31), <u>shoda</u> * (0.28)
manual IPL * [14] <u>IPL</u> (0.54), <u>manuální</u> (0.46)	matched [11] <u>odpovídat</u> (0.23), <u>za</u> (0.23), <u>nalezený</u> (0.18), <u>další příkazy</u> * (0.12), <u>splňovat</u> (0.12), <u>vyhovět</u> (0.12)
manual mode * [26] <u>režim manual</u> * (1.00)	matches [56] <u>odpovídat</u> (0.85), <u>souhlasit</u> (0.15)
manufacturer * [10] <u>výrobce</u> * (0.83), <u>zařízení IBM</u> * (0.17)	matching * [13] <u>odpovídající</u> * (0.63), <u>odpovídat</u> (0.37)
many [404] <u>mnoho</u> (0.87), <u>kolik</u> (0.13)	material * [30] <u>materiál</u> * (1.00)
map [12] <u>mapovat</u> (0.51), <u>AS/400</u> (0.28), <u>datové typy</u> * (0.21)	materials * [11] <u>materiály</u> * (0.60), <u>materiál</u> * (0.40)
map * [31] <u>mapa</u> * (0.68), <u>map</u> * (0.32)	matrix [16] <u>matice</u> * (1.00)
	max [41] <u>max</u> (0.79), <u>maximálně</u> (0.21)
	maximum * [137] <u>maximum</u> * (0.52), <u>maximálně</u> (0.48)
	maximum length * [18] <u>maximální délka</u> * (0.72), <u>maximální délka parametru</u> * (0.28)

Figure 3: Sample of English–Czech computer oriented dictionary extracted from a parallel corpus.

rachot * [11] crashing (0.61), its (0.39)	razit [14] quickly (0.55), wrong (0.45)
rada * [33] never (0.44), make (0.30), kids* (0.27)	rád [345] <u>love</u> (0.66), <u>loved</u> (0.19), don't (0.15)
raději [136] <u>rather</u> (0.53), <u>better</u> (0.29), <u>prefer</u> (0.18)	rádus * [21] radio* (0.71), disc* (0.29)
radio * [18] <u>radio</u> * (1.00)	rámec * [11] companies* (0.55), endometrial (0.45)
radit [82] <u>told</u> (0.67), how (0.33)	rána * [53] hole* (1.00)
radnice * [13] <u>city hall</u> * (0.55), day* (0.45)	ráno * [57] <u>morning</u> * (0.69), wound* (0.16), girls* (0.15)
radost * [67] <u>joy</u> * (0.74), <u>happiness</u> * (0.26)	rány * [19] wounds* (0.53), blows* (0.47)
radovat [24] joy* (0.53), well (0.47)	readers [10] german (0.34), again (0.22), american (0.22), those (0.22)
radý * [26] current article text* (0.58), <u>advices</u> * (0.42)	reader [20] digest (0.39), <u>reader's</u> (0.32), year* (0.15), nearly (0.15)
rajčata * [11] <u>tomatoes</u> * (0.57), tomato (0.43)	reagovat [108] <u>respond</u> (0.69), <u>react</u> (0.31)
raketa * [10] <u>rocket</u> * (1.00)	reakce * [26] <u>reaction</u> * (0.71), something* (0.29)
raketoplán * [11] <u>shuttle</u> * (1.00)	recept * [13] <u>recipe</u> * (0.57), <u>prescription</u> * (0.43)
rakety * [10] <u>missiles</u> * (0.59), missile* (0.41)	reeves [13] <u>reeves</u> (0.82), scott's (0.18)
rakev * [10] husband* (0.54), <u>coffin</u> * (0.46)	reid [11] <u>reid</u> (0.84), golf balls* (0.16)
rakovina * [96] <u>cancer</u> * (1.00)	reklama * [23] <u>advertising</u> * (1.00)
rakovina plic * [16] <u>lung cancer</u> * (1.00)	rentgen * [10] rays* (0.50), x (0.50)
rakovina prsu * [22] <u>breast cancer</u> * (0.74), cancer* (0.26)	republikáni * [11] gingrich (0.60), showing (0.20), divorce* (0.20)
ralph [41] <u>ralph</u> (0.81), jaymee (0.19)	respektovat [10] respect* (0.53), got (0.47)
ramena * [59] <u>shoulders</u> * (0.81), arm* (0.19)	respirátor * [12] <u>respirator</u> * (1.00)
rameno * [70] <u>shoulder</u> (0.83), right (0.17)	restaurace * [51] <u>restaurant</u> * (1.00)
ranc * [22] <u>ranch</u> * (0.75), years* (0.25)	rezervace * [16] <u>park</u> * (0.67), elephants* (0.33)
raul [17] <u>raul</u> (0.85), dad* (0.15)	režisér * [10] <u>spielberg</u> (0.59), <u>director</u> * (0.41)
ravussin [10] fat calories* (0.26), she's (0.26), sanchez (0.25), carbohydrates (0.23)	

Figure 4: Sample of Czech–English fiction dictionary extracted from a parallel corpus.

References

- BROWN, PETER F.; DELLAPIETRA, S. A.; DELLAPIETRA, V. J.; MERCER, ROBERT L. (1993), "The Mathematics of Statistical Machine Translation: Parameter Estimation". In *Computational Linguistic*, 19(2): 263–331.
- CUŘÍN, JAN. (1998), "Automatická extrakce překladu odborné terminologie." *MSc Thesis*, Institute of Formal and Applied Linguistic, Charles University, Prague. 89 pp. (in Czech)
- ČMEJREK, MARTIN. (1998), "Automatická extrakce dvojjazyčného pravděpodobnostního slovníku z paralelních textů". *MSc Thesis*, Institute of Formal and Applied Linguistic, Charles University, Prague. 82 pp. (in Czech)
- GALE, WILIAM A.; CHURCH, KENNETH W. (1993), "A Program for Aligning Sentences in Bilingual Corpora". In *Computations Linguistic*, 19(1): 75–102.
- HAIČ, JAN; HLADKÁ, BARBORA. (1998), "Tagging Inflective Languages: Prediction of Morphological Categories for a Rich, Structured Tagset". In *Proceedings of Coling/ACL'98*, Montreal, Canada.
- WU, DEKAI; XIA, XUANYIN. (1994), "Learning an English–Chinese Lexicon from a Parallel Corpus". *Association for Machine Translation in the Americas*, Oct. 94: 206–213, Columbia, USA.