

Language Technology Tools in the Translator's Practice

Gábor Prószéky

MorphoLogic, Budapest, Hungary

An open, extendible multi-dictionary system is introduced in the paper. Simultaneously an unlimited number of dictionaries can be held open, thus by a single query step, all the dictionaries (translations, explanations, synonyms, etc.) can be looked up. The implemented system (called MoBiDic) knows morphological rules of the dictionaries' languages. Thus, never the actual (inflected) words, but always their lemmas - that is, the right dictionary entries - are looked up. MoBiDic has an open, multimedia architecture. It has been designed for translator workgroups using intranet or internet.

Keywords: language technology, machine-aided translation

1. Introduction

It is almost a commonplace that texts - books, newspapers, letters, official memos, brochures, any type of publications, reports, etc. - in the nineties are written, sent, read and translated with the help of the electronic media. Consequently, traditional information sources, like paper-based dictionaries, and lexicons, are no longer as much a part of the translation environment. Electronic dictionaries for most developers just mean, however, to make the well-known paper dictionary image appear on the computer screen. It is easy to understand why we say that dictionary computerization does not mean producing machine-readable versions of traditional printed dictionaries, but the combination of the existing lexical resources with up-to-date language technology. On the other hand, there is a question whether we have to continue in the traditional way of developing new - and different - lexicons for any new application/system, starting from scratch every time and therefore consuming time, money and manpower, or is it

new lexicons. In what follows, timely to think of the possibility of making the effort to converge, trying to avoid unnecessary duplications and - where possible - building on what already exists (Calzolari 1994). Consequently, in the near future we have to combine the two above needs: making existing lexical resources computationally accessible and showing the strategy how to develop we try to argue for changes in development strategies of electronic translation dictionaries. Today's language technology can - and must - use dynamic actions, like morpho-syntactic analysis, lemmatization, spell checking, and so on. On the other hand, dictionaries can never be full in any sense, therefore we have to make parallel multidictionary access possible. It means that a single dictionary look-up should use an unlimited number of lexical resources that are available for the translator.

2. An intelligent dictionary look-up system

To start with the most natural activity concerning dictionaries is searching them for a single word. There is no problem if it can be found among the headwords of the dictionary, that is, when the input string can match. But sometimes the translator starts the look-up process by clicking an inflected word-form of an open document that cannot be found among the headwords. For the user it is a boring and time-consuming task to type the lexical form, that is, the one accepted letter-by-letter by the dictionary. To make the system able to find the stem of the input word-form automatically, MoBiDic uses a lemmatizer that provides the dictionary

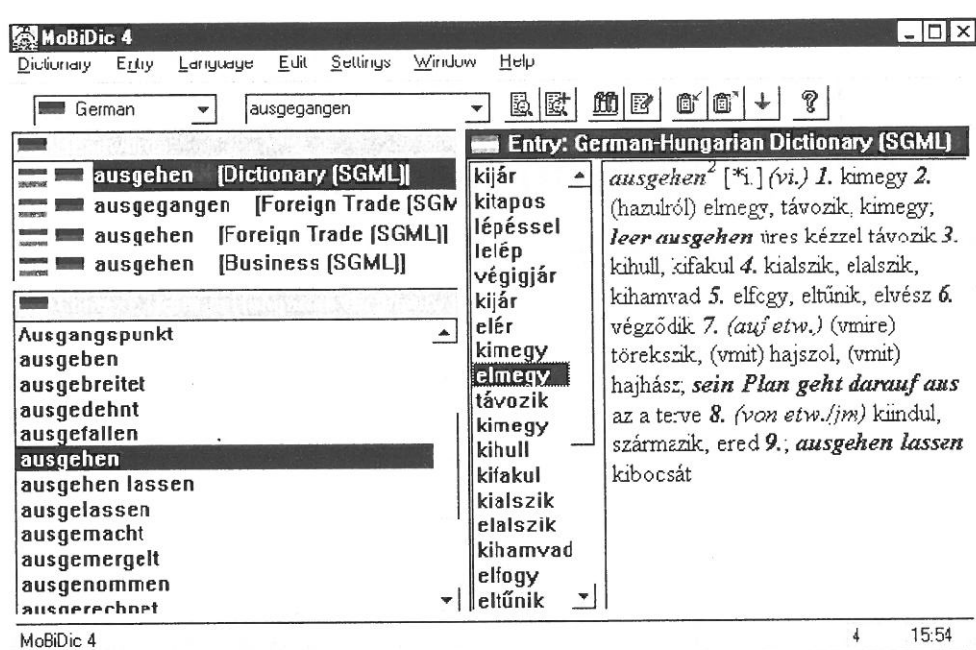


Fig. 1. Look-up of a morphologically complex inflected form: 'ausgegangen' in a German-Hungarian dictionary.

look-up module with the stem(s) to be found (Fig. 1).

Translators frequently want to find the word as a part of *multi-word expressions* or idioms. If the user does not know whether the actual word is part of some phrasal compound or idiom, the traditional paper dictionaries are very difficult to use. Namely, if the word in question is the so-called headword of a multi-word expression, it can be found easily. In case it is not the headword, one has to know the phrasal compound the word is a part of, but it is a typical "Catch 22" situation: if the expression is known why to search the dictionary for it? MoBiDic helps the user to find all the multi-word expressions containing the actual word's stem, independently whether it is a headword or not. E.g. not only 'lead' but both 'dog' and 'life' provide us (among others) with the multi-word expression 'lead a dog's life' that can be found under 'lead' only in a paper dictionary. In other words, users of the traditional dictionaries are supposed to know the expression (what's more: the keyword of the expression) in order to find it in the lexicon. Search for 'lead a dog's life' through its components gives the following result in MoBiDic:

lead {*lead, leads, leading, led*}

27 occurrences in expressions of the basic dictionary,

dog {*dog, dogs, dog's, dogs'*}

21 occurrences in expressions of the basic dictionary,

life {*life, lives, life's, lives'*}

77 occurrences in expressions of the basic dictionary,

lead AND life

5 occurrences in expressions of the basic dictionary,

dog AND life

2 occurrences in expressions of the basic dictionary,

lead AND dog

1 occurrence in expressions of the basic dictionary,

lead a dog's life

1 occurrence as an expression in the basic dictionary.

'Bi' is somewhat misleading in the name MoBiDic. Bilingual in this sense means that the source and the target language are not the same types of object for the program. For MoBiDic, source language is the language the *morphology* of which has to be known, to provide the user with adequate output. The output is expected to be in the target language - the characters, the alphabetic order, etc. of which has to be known to make the hits appear on the screen in adequate format. Of course, the source and target

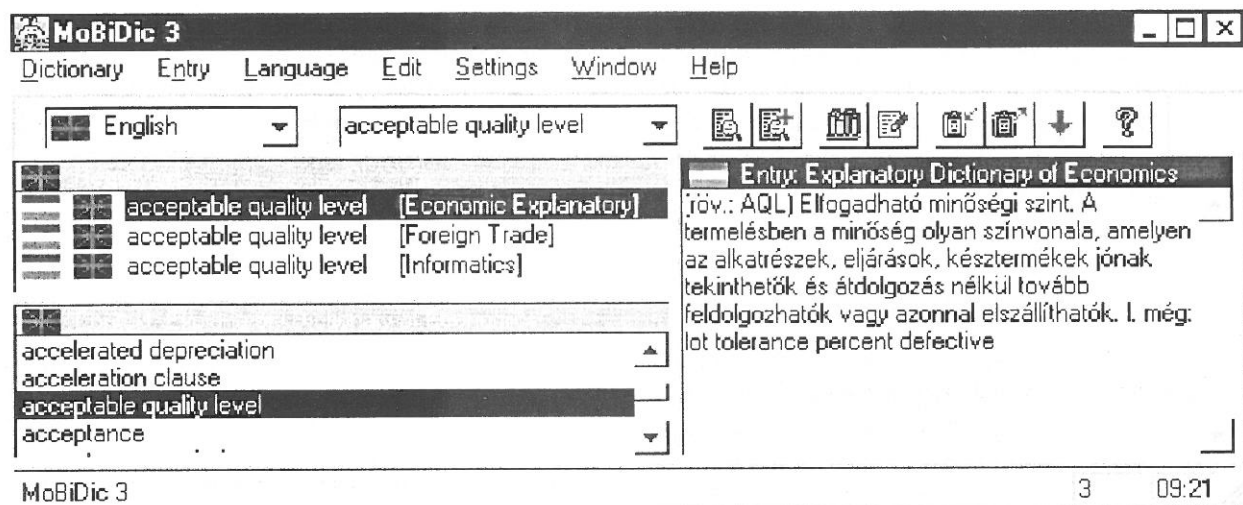


Fig. 2. Hungarian explanation of 'acceptable quality level' in the *English-Hungarian Economical Explanatory Dictionary*.

languages can be the same, e.g. in *explanatory or etymological dictionaries* (Fig. 2).

There is another sort of monolingual dictionary, the synonym dictionary. The translator frequently wants to use a synonym (antonym, hypernym, hyponym) of the actual word. An intelligent software tool, like MorphoLogic's Helyette¹, is the combination of a thesaurus (synonym dictionary), a morphological analyzer and a generator, because the output is

re-inflected according to the morphological information contained by the input word-form. The - so-called inflectional - thesaurus works as follows:

INPUT:	came
ANALYSIS:	came = come + Past
STEM:	come
SYNONYM:	go
SYNTHESIS:	go + Past = went
OUTPUT:	went

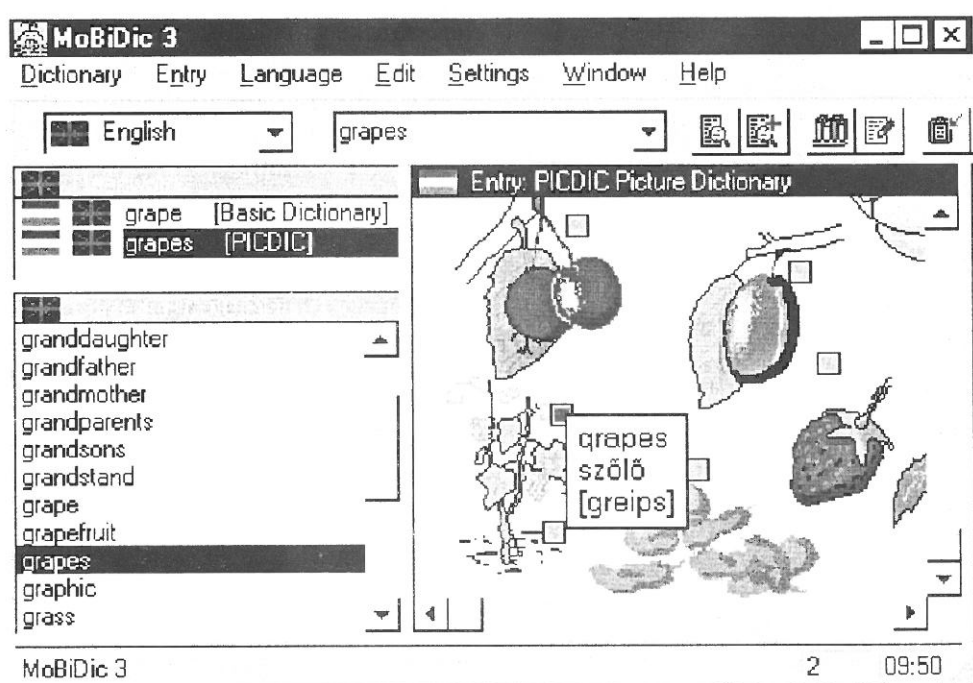


Fig. 3. 'grapes' (from the PicDIC picture dictionary) with pronunciation in MoBiDic.

¹ To be combined with MoBiDic in the near future.

There are special sorts of information in a dictionary. For example, pronunciation is not typically needed for translation, but can be useful for language learners. Pronunciation of the word is, therefore, an information that should be switched on and off, according to the user's needs. In an electronic dictionary it is expected that not only the written phonetic transcription, but also the *spoken* output can be heard. If the dictionary supports multimedia, explanatory *pictures* can help understand the word, even for professionals, not for language learners only (Fig. 3). If the translator makes a spelling error, first a *speller* starts, and then the corrected word-form is sent to the dictionary lookup system. In large, professional paper dictionaries examples usually belong to the entries. In electronic dictionaries, occurrences of the same word in real contexts — sometimes with translations — are easy to show with the help of available monolingual or aligned bilingual corpora. In order to search corpora like dictionaries, texts and lexicons should be represented in a uniform way: with the help of SGML, that is, the Standard Generalized Markup Language.

3. Dictionaries in SGML

The lexicographic basis for MoBiDic is supplied by various publishing houses. More precisely, MorphoLogic has licenses to almost 50 dictionaries already published in paper format of miscellaneous topics, diverse sizes and many language pairs. The user can choose which dictionary to use in general and which of them open actually. Currently, if all the available dictionaries are open, MoBiDic handles approximately 1 million lexical entries. Some of the dictionaries, mainly the terminological ones, have usually a very simple list-based structure. Dictionaries shown by Fig. 1 and Fig. 2, however, appear on the screen with the traditional paper dictionary image. It is done by using SGML representations and an on-line SGML-RTF conversion. MoBiDic can do exact structural search not influenced by the layout at all. The original lexical resource - even it has been available in electronic format - did not use SGML. A special system for (semi-)automatic conversion of some formatted text files containing dictionary

data to SGML format has also been developed for the MoBiDic environment. This utility system is not available for the end-users, it serves industrial purposes². First, in order to enable selective access to the information in dictionary entries, a thorough structural analysis is done, while inconsistent and faulty entries are marked. They are corrected later, manually. The resulting SGML-annotated dictionaries (Fig. 5) are enhanced with the necessary indexes. They are lemma-variants and expanded sub-entries made with the help of existing language technology modules (Prószéky 1994).

Users like to work with their own little vocabularies, glossaries, and the professional translator is usually asked to use official translation equivalents provided by the employer. These glossaries are generally never published, but there is a need to use them in the same environment. MoBiDic is able to treat user dictionaries containing any type of information sources (lexicons, encyclopedias and dictionaries).

The strength of this method is that user dictionaries are looked up for a word exactly when other dictionaries, thus translator's remarks can also be read when other dictionaries provide the user with their translation equivalents. Here we have to emphasize again that MoBiDic is not yet another electronic dictionary, but a multi-dictionary environment where a single word is sent to every open dictionary by a single mouse-click. In Fig. 4 the user started from the word-form '*duties*', and eight dictionaries (that are open and contain English on either the source or the target side) send translations to the screen.

4. Implementation features

The most recent development is MoBiDic's client-server implementation. Its server side (WinNT, Unix and Novell) consists, in fact, of two servers: the linguistic server and the dictionary server. The user interface and screen handling modules take place on the (Win32, Unix/Linux, etc.) client side. There are many software modules of other vendors on the market that can also be combined with MoBiDic through its well-defined *application programming interface* (API). With the help of this API,

² See http://www.morphologic.hu/e_sgml.htm

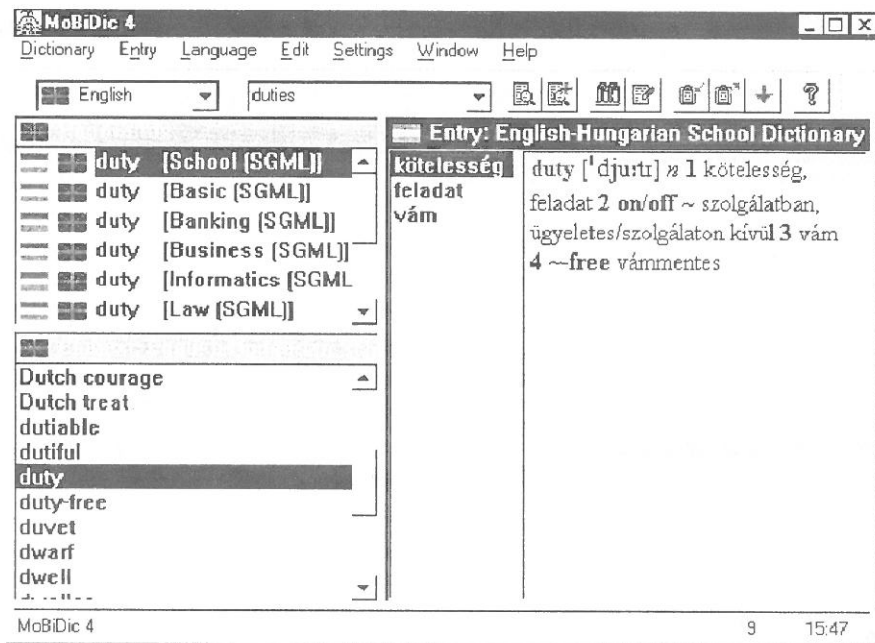


Fig. 4. Search for the (lemma of) 'duties' in a set of English-Hungarian dictionaries.

the user can communicate to the other modules from MoBiDic without leaving it. Because of technical and legal reasons, it can, of course, be done in collaboration with the developer of the product in question. The picture dictionary shown by Fig. 4 is a working example: the vocabulary part of the (also commercial) CALL program called PicDIC is available for MoBiDic users from the familiar environment. Translators who generally use their favorite word-processor while translating can use MoBiDic from their word-processing tools with the help of the included macros. Another important issue is that users can use their CD-ROM drive for other purposes while translating. Namely, MoBiDic has minimal space requirement because of its compression method³, therefore the full dictionary system can be copied to the hard disk: thus the CD drive is freed and can be used for other purposes.

5. Comparison with other methods

There are several dictionary programs both in laboratories and on the market, but only some of them share the so-called "intelligent" features with MoBiDic. Rank Xerox developed in the COMPASS and LOCOLEX projects a

prototype that accesses enhanced and structurally elaborated dictionaries with an intelligent, context-sensitive look-up procedure, presenting the information to the user through an attractive graphical interface. (Feldweg and Breidt 1996) Unlike MoBiDic, it does not have access to more than one dictionary at the same time. Consequently, user dictionaries are not supported. SGML is, however, used both in the dictionary and in the corpus modules. There is a focus on the intelligent treatment of multi-word units in the IDAREX formalism (Breidt et al 1996). Another project with similar aims is GLOSSER. Its prototype (Nerbonne et al. 1997) carries out a morphological analysis of the sentence in which the selected word occurs and a stochastic disambiguation of the word class information. This information is then matched against a (single, but SGML) dictionary and corpora. The GLOSSER prototype displays context dependent translations and on request, examples from the available corpora. Neither of the above developments nor other web dictionary services (e.g. WordBot) shares all the important features with MoBiDic: client-server architecture, multi-dictionary access, user dictionary handling, parallel (and intelligent) dictionary and corpus look-up. What's more, MoBiDic is commercially also available, that is tested by thousands of "real" end-users.

³ Average 1-2 MB/dictionary.

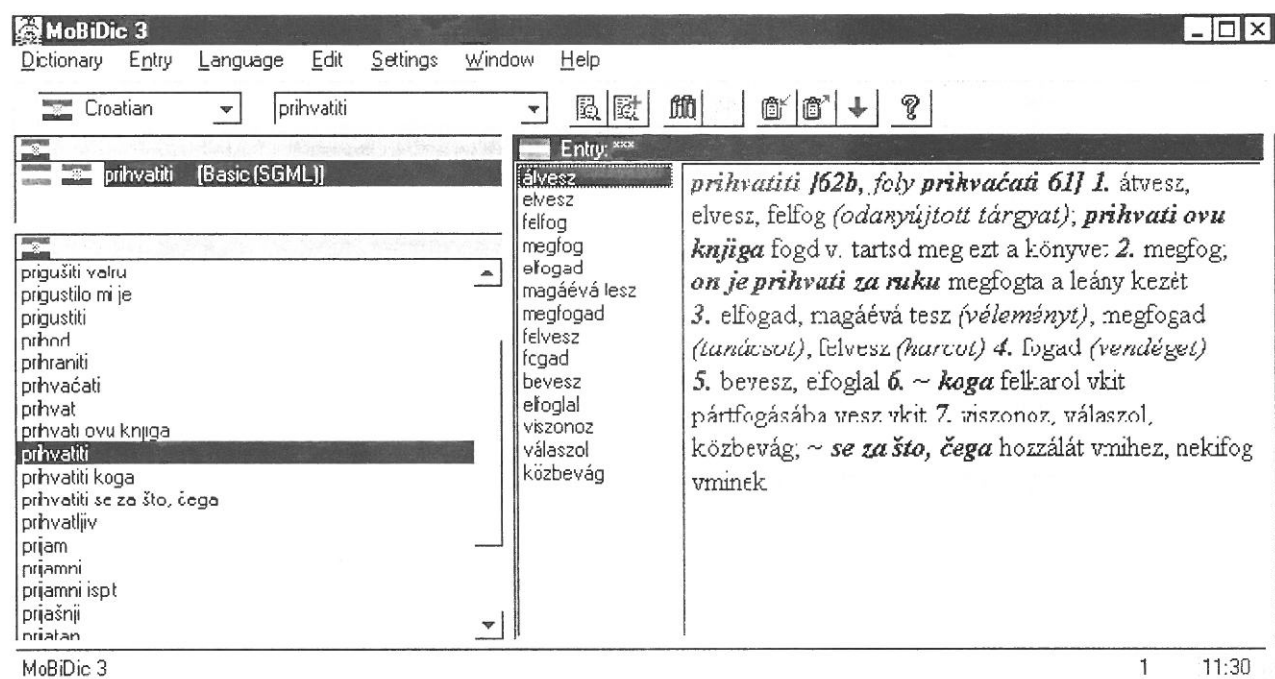


Fig. 5. A Croatian-Hungarian Dictionary in SGML.

6. Conclusion

MoBiDic is a multi-dictionary translation environment based on a client-server architecture. It consists of the following main parts: linguistic server, dictionary server and the client with the graphical user interface. The system has several features that makes it different from other translation tools:

- the linguistic server is dictionary-independent and language-dependent⁴;
- the dictionary server has intelligent access to various sorts of dictionaries (from SGML to multimedia) and bilingual corpora;
- simultaneously an unlimited number of dictionaries can be held open, thus by a single interrogation step, all the dictionaries (with translations, explanations, synonyms, etc.) can be surveyed;
- the translators' own glossaries built with the help of the system may also be disseminated (as new dictionaries, with the needed copyrights) among other users, if needed;
- it has an open architecture and a well-defined API;

- it has been implemented and is available with a gradually increasing number of dictionaries for numerous language pairs.

MoBiDic is, consequently, not a research project only, but an attempt to provide both the professional translators and the language learners with a set of practical translation tools.

References

- [1] BREIDT E., F. SEGOND AND G. VALETTO, (1994), "Local Grammars for the Description of Multi-Word Lexemes and Their Automatic Recognition: in Texts", *Complex-94 - Papers in Computational Lexicography*, Linguistics Institute, HAS, Budapest, pp. 19-28.
- [2] CALZOLARI, N., (1994), "Issues for Lexicon Building", In: A. Zampolli, N. Calzolari & M. Palmer (eds.) *Current Issues in Computational Linguistics: In Honour of Don Walker*, Kluwer / Giardini Editori, Pisa, pp. 267-281.
- [3] FELDWEIG, H. AND E. BREIDT, (1996), "COMPASS - An Intelligent Dictionary System for Reading Text in a Foreign Language". *Complex-96 - Papers in Computational Lexicography*, Linguistics Institute, HAS, Budapest, pp. 53-62.
- [4] HUTCHINS, J., (1996), "Introduction", *Proceedings of the EAMT Machine Translation Workshop*, Vienna, pp.7-8.

⁴ Recently, English, German, Hungarian, Polish, Czech and Romanian morphological components are available for the MoBiDic users. Descriptions for further languages are under development, see <http://www.morphologic.hu> for the actual list of languages.

- [5] KINGSCOTT, G., (1993), "Applications of Machine Translation". In: Kohn, J. (ed) *Transferre necesse est... (Current Issues of Translation Theory)*, Szombathely, pp. 239–248.
- [6] NERBONNE, L. KARTTUNEN, E. PASKALEVA, G. PRÓSZÉKY AND T. ROOSMAA, (1997), "Reading More into Foreign Languages", *Proceedings of the Fifth Conference on Applied Natural Language Processing*, Washington.
- [7] PRÓSZÉKY, G., (1994), "Industrial Applications of Unification Morphology", *Proceedings of the 4th Conference on Applied Natural Language Processing*, Stuttgart, pp. 157–159.

Received: March, 1999

Accepted: June, 1999

Contact address:

Gábor Prószéky

MorphoLogic

Késmárki u. 8.

H-1118 Budapest

Hungary

e-mail: proszeky@morphologic.hu

GÁBOR PRÓSZÉKY (1954) has backgrounds both in computer science and linguistics. He holds a Ph.D. in computational linguistics. He has written more than 50 scientific publications mainly on human language technologies. Since 1991 he has been the director of MORPHOLOGIC, the first and only language industrial company in Hungary (see <http://www.morphologic.hu>).
