# PERSONAL TRANSLATION APPLICATIONS

Ian Johnson
Sharp Laboratories of Europe, Oxford, UK
ianj@sharp.co.uk

Two new language technology products from Sharp Laboratories of Europe in Oxford are presented – a new English-Japanese dictionary product which has recently been released in Japan on Sharp's Shoin word processor and a multilingual document generator application which is currently available on the Sharp Font Writer. These products address two problems in language technology: the browsing of documents in English on the Internet when you are a native speaker of another language and the generation of standard documents in a language that you may not know at all. These are the first products of the Language Technology group at SLE in Oxford and we have plans to improve and extend them in a variety of ways.

## BACKGROUND

Sharp's interest in language technology products dates back to the early 1980s when the DUET E/J family of translation products was started. The first product sold in 1987 was based on a Unix workstation and included an OCR. The cost of the system was in the region of £20k and the target market was professional translators and translation agencies. The system has been refined over the years – in particular it has been integrated with a Web browser – and the latest Power E/J system is now available with Sharp's PC or as a standalone software package for Windows 95 selling at around £60. The increasing affordability of PCs has opened up the market for these translation products and the advent of the world wide web has led to an increased demand for on-line translation software which can be used by people who may not be professional translators but who would like to access information in other languages. English is the dominant language on the web and for that reason Sharp are focussing on browsing and translation applications from English into the target languages. Our first language pair is English-Japanese but further language pairs are planned for the near future.

Sharp Laboratories of Europe Ltd (SLE) based in Oxford was established in 1990 as Sharp's first R&D laboratory outside Japan. Our first product was the Multilingual Document Generator (MDG) which is an application supplied with the Sharp Font Writer 750/760. The objective behind Sharp's development of the Multilingual Document Generator (MDG) was to provide users who have little or no knowledge of foreign languages with an application which would allow them to customise standard documents written in their own language and to automatically generate translations of these documents in other languages with guaranteed linguistic, stylistic and cultural correctness. This would enable for example small businesses and individuals in the UK  to overcome language barriers encountered either in their business

or personal lives and allow them to communicate with people who may not necessarily understand English. The MDG provides them with the ability to communicate information about their business or their personal lives in the language of the recipient which will improve understanding and create a good impression, both important factors in an increasingly competitive global trading community. The MDG output needed to be 100% correct, as the users would not normally be able to assess the quality of the text generated. The first version of the MDG on the Sharp Font Writer 750 was developed by Sharp Laboratories of Europe Ltd based in Oxford in collaboration with Sharp Manufacturing Company of UK Ltd, based in Wrexham. The latest UK version of the system on the Font Writer 760 generates documents prepared in English into French, German, Spanish and Italian.

In providing these products, Sharp is attacking the problem of computerised translation of documents from two perspectives. In the case of the application allowing users to browse English documents and see words and phrases they may not know translated in their own language, the aim is not to achieve 100% correct translation, merely to allow users to combine their own knowledge of the language with the information provided by the system to gain an understanding of the text concerned. In the case of the MDG, there is an absolute requirement to ensure that the document is 100% correct because it is assumed that users themselves may not be able to judge the accuracy of the translations generated. With translation into the user's own language we can rely to a large extent on the user being able to piece together the words and phrases presented by the system into a coherent whole. So we have provided the user with two translation tools which within certain constraints perform well. Future work will be aimed at improving the accuracy of the browsing tool so that it comes closer to full translation and providing the user of the MDG with more flexibility in adapting the template documents to individual requirements whilst at the same time maintaining 100% accuracy in generating the translation.

## MULTILINGUAL DOCUMENT GENERATOR (MDG)

The idea for the Multilingual Document Generator dated back to an early research paper by Tomita (Saito and Tomita, 1985). A possible product was contemplated in 1989 when the author was working at Lexpertise S.A., a spinoff company of Automated Language Processing Systems (ALPS). Sharp produced an application allowing Japanese people to generate letters in English for the Shoin word processor in 1989. A demo of a prototype in 1991 at UMIST funded by BT and implemented by Jones & Sager (see their earlier technical report dated 1989) gave further stimulus to the idea. Since then further products have begun to appear on the market. Examples of systems which we evaluated were LinguaWrite (Multilingua 1989-1994) and Ambassador (1992, Catena Corporation/Language Engineering Corporation) and the 1996 version of Sharp's own system on the Shoin word processor. Many books have been written about letter writing including bilingual letter writing. Although we did not make use of the information in them, the following books provided a particularly interesting source of ideas for the product: Ferney et al (1983, 1990), Harvard et al (1989, etc), Davies et al (1989).

The MDG application described here was developed for the Sharp Font Writer, a compact, portable word processor with full PC-sized LCD screen and integral, near laser quality printer. Major constraining factors in developing the software were the lack of a hard

disk and the relatively slow processor. The Multilingual Document Generator allows users who may not be able to write a document in a foreign language themselves to adapt standard template letters in their own language, select the desired target language from the list of languages supplied with the system and generate a translation of the document in the selected language which is guaranteed to be linguistically, stylistically and culturally correct. Approximately 60 standard document templates are supplied with the system which currently handles English, French, German, Spanish and Italian. These 60 templates allow users to make many modifications to reflect their own requirements, thus providing them with thousands of variations on these standard documents, all with guaranteed correct translations. The templates currently provided with the system are shown in Figure 1 below:

```
MULTILINGUAL DOCUMENT GENERATOR        Personal
        Document Types                         Announce birth of baby
                                               Congratulate someone on promotion or new
  Meetings                                     job
        Propose a meeting                      Congratulate someone on new company or
        Confirm a meeting                      office
        Change a meeting                       Congratulate someone on retirement
        Cancel or postpone a meeting           Congratulate someone on personal event
        Send arrival information               Express condolences to business colleague
        Agree to meet visitor on arrival       Express condolences to personal friend
  Hotel Reservations                           Thank someone for hospitality
        Request basic hotel information        Thank someone for gift
        Request detailed hotel information  Invitations
        Make hotel reservation                 Invitation to personal events
        Confirm hotel reservation              Invitation to public events
        Modify hotel reservation               Invitation to wedding ceremony
        Cancel hotel reservation               Invitation to wedding reception
  Orders                                       Acceptance of invitation
        Place order                            Rejection of invitation
        Acknowledge receipt of order        Addresses and Contacts
        Cancel order                           Change of address
        Cancel delayed order (supplier's       Change of job and address
        fault)                                 Retirement and change of address
        Cancel delayed order (not supplier's   Announcement of new subsidiary
        fault)                                 Announcement of new company
        Complain about goods received       Applications for Jobs and Courses
  Payments                                     Curriculum vitae
        First reminder for payment             Request information about course
        Payment already sent                   Application in response to job advert
        Apologise for delayed payment          Unsolicited job application
        Acknowledge receipt of payment         Acknowledge receipt of application
                                               Request further information from applicant
  Product Information                          Invite to interview
        Request product information            Make offer
        Send product information               Reject application
                                               Accept offer
                                               Decline offer
```
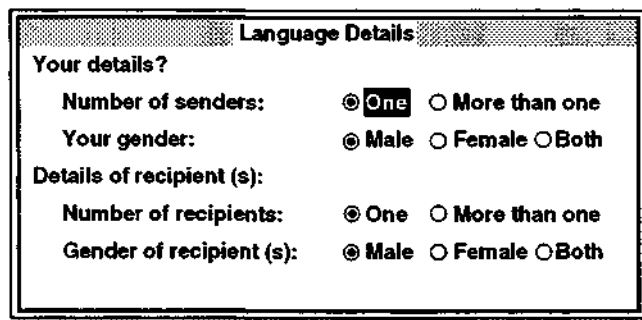
Figure 1: List of Document Templates Provided with the MDG


In designing the MDG it was quickly apparent that the user interface was very important. Some products provided flexibility by making the user interface very busy which made the text difficult to read. We therefore decided for a first version to make the system as simple as possible to use, bearing in mind that the target market for the Font Writer is composed largely of people who are not very familiar with computers or word processors. This meant that instead of opting for our original idea of constructing a document from individual sentences and paragraphs, we decided instead to design full document templates so that we could maintain more control over the flexibility offered to users and keep the system simple and straightforward to operate. A further benefit of this approach was that it was possible to maintain the overall structure and flow of the documents. It would have been very difficult to achieve this by allowing users to select sentences freely from a database, because the need to insert link words and phrases (e.g. 'first of all', 'however', 'finally') and to

maintain a consistent tone and style would have made it impossible in all cases to guarantee 100% correct translation. The disadvantage of this approach, however, is that the system is currently quite restrictive in the degree to which users can modify the templates provided.

Existing Machine Translation technology is not reliable enough for documents to be translated automatically with a guarantee that the translations are completely correct. In applications involving translation into one's own language that need not be a problem because generally people are able to piece together the meaning of a document even if the translation is not entirely accurate. However, in generating documents from one's own language into another language, users cannot check on the quality of the output themselves. This fact determined the overall design of the MDG application and distinguishes it from full Machine Translation. We do not allow the user to enter their own free text, but require them instead to select from a database of standard template documents. Within this restriction the user is presented with various possibilities of modifying the document in a controlled way, thus maintaining the guarantee of 100% correct output. Other existing products such as LinguaWrite enable users to construct documents from words, phrases and sentence fragments that have pre-stored translations. However, the user has to have an excellent command of the foreign language in order to join these segments together in a grammatically correct and stylistically accurate way. The MDG application manages to provide a user who may have little or no knowledge of the foreign language with a fair degree of flexibility in adapting the standard document templates, whilst maintaining the requirement of generating accurate translations.

The system works as follows. Having selected a document category from the list presented to the user when the system starts up and then selected the specific document required, the user enters details of his/her address which is saved as a default address for future documents. Up to three sets of personal details (work, home, other) can be stored in this way. Once the user has entered or selected the address required, the system presents the language details dialogue box shown in Figure 2:



Figure 2: Language Details Dialogue Box

In order to generate the appropriate agreement in the translated text, it is necessary to know the number and gender of the sender and recipient. The MDG then goes into Document Creation Mode which allows the user to view the document selected and edit certain parts of it to suit his/her needs. The sample document created by the MDG shown in Figure 3 contains three types of fields or portions of text: fixed text field, details text field and editable text

field.



John Smith
Smith Engineering
Sharp Street
Newtown
Borchester
B99 9MM
England
Telephone:     0987 654321
Fax:           0987 654322
Mobile Phone:  0876 543210
Email:         0765 432109

**Details Text Fields**

Details Text is text that you entered in the Your Personal Details Dialogue box. This contains information obtained from your Home, Work or Other Details.
If you press the [ Edit ] key or the [Return] key when the cursor is on one of these fields, the Details Dialogue box will be displayed and the cursor will be on the corresponding field.
You can also alter any other fields contained in this box by using the [ Tab ] or the [ Shift ] + [ Tab ] keys to move to the relevant fields.
When the [Return] key is pressed, the Details Text fields will be updated.

<FIRST NAME> <LAST NAME>
<JOB TITLE>
<ORGANISATION>
<HOUSE NUMBER AND STREET NAME>
<PLACE NAME> <POSTCODE>
<REGION>
<COUNTRY>

— **Recipients Details**

**Editable Text Fields**

You can change an Editable Text field to suit yourself. The different types of Editable Text field are discussed in the Operation Manual. You can only edit one field at a time.

Our reference: <REFERENCE>
Your reference: <REFERENCE>

<TODAY'S DATE>

Dear Sir,

I will be visiting your area from Monday <DATE> until Friday <DATE> and would like to meet you for a brief discussion at some time during this period <DATE>. Could you let me know if you could fit such a meeting into your schedule? <___> <___>

I look forward to hearing from you.

Yours faithfully

**Fixed Text Fields**

Fixed Text fields are automatically created by the MDG and cannot be edited by the user.

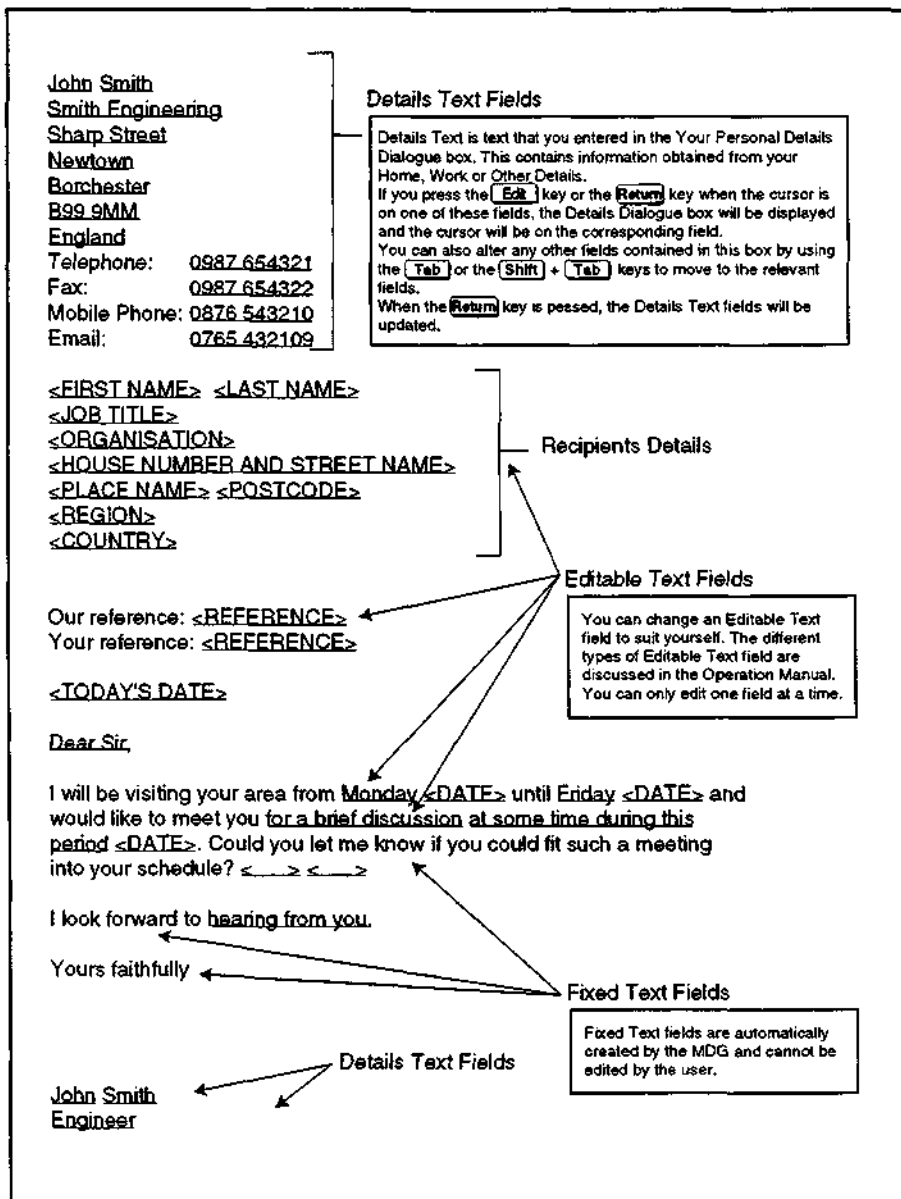**Details Text Fields**

John Smith
Engineer

Figure 3: Sample Document from Font Writer 760 Manual

The fixed portion of the text cannot be altered by the user. The personal details field are entered automatically into the document once the user has registered his/her default address details with the MDG system. The editable text field allows users to enter their own text such as dates, numbers, names, etc. This text is transferred without alteration into the target language.

The flexibility of the system is achieved principally by means of option slots which permit the user to select one or more items from a list of possibilities and, where appropriate, to enter their own text. Three examples, illustrating single choice slots, multiple choice slots, and multiple choice slots with optional user input are shown below:
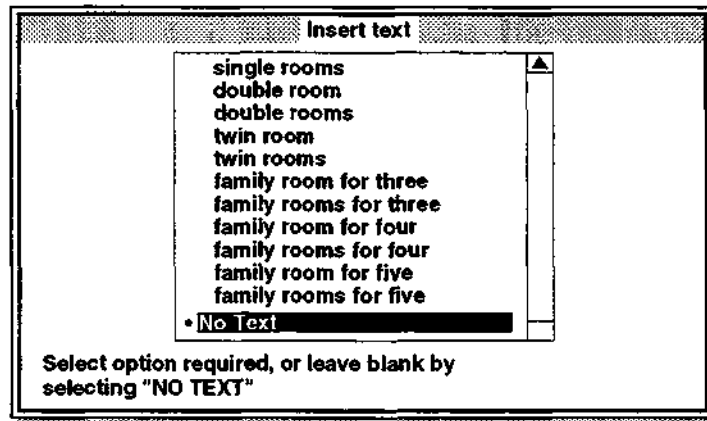
Figure 6: Multiple Option Slot with User Input
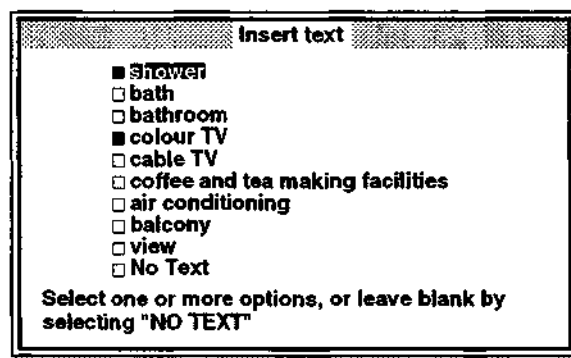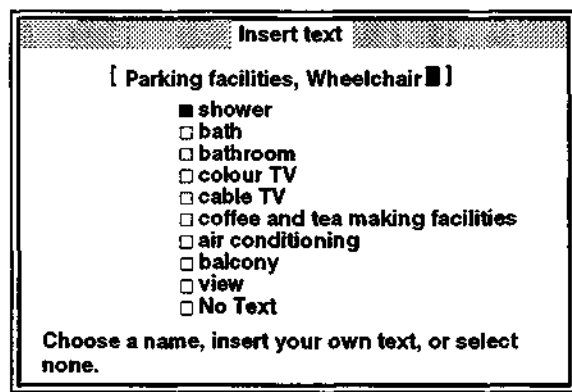
Figure 4: Single Option Slot



Figure 5: Multiple Option Slot



Some sentences are optional and can be omitted if required. Once the user has entered the data required in the editable fields, selected the appropriate items in the various option slots and decided whether or not to include any optional sentences, the document is ready to be generated in the target language. The user simply selects the language required from the list of languages available (currently English, French, German, Spanish and Italian) and the translated version of the document is generated automatically in word processor format, allowing the user to reformat or make any further changes before sending the document to the recipient (see Appendix for sample output).  The document can be saved for future reference

or reuse.

We encountered a number of translation problems by adopting this approach. The major constraint on the data side was that the languages needed to be tied very closely together. The advantage of doing this was that the data could be used multi-directionally, i.e. the same Spanish data could be used for English-Spanish and Spanish-German, etc. This meant, however, that the number of each type of field had to match up exactly in each language and that each element of a list of options had to have an equivalent. This was not always possible and in fact this problem gets worse if you carry on adding languages to the system. Future versions of the system will therefore have to find a way of loosening this link between the languages to take account of this. One example of this is demonstrated by the need for one-to-many equivalences to handle the equivalence between "Frau" and "Mrs'/"Ms". Another problem was with prepositions and determiners before user entered place names. An example of this is given below, showing the need to change the preposition and determiner depending on the number and gender of the place name:

---

We would like to invite you to {the conference/the exhibition/0} [TITLE] to be held at [PLACE]

... die im [Queen's Conference Centre] stattfinden wird.
... die in [Queen's Conference Centre] stattfinden wird.

... che si terrà allo/alla/all'/agli/alle [PLACE].
... che si terrà in nel/nella/nell'/negli/nelle [PLACE]

---

Figure 7: Example of problems with prepositions and determiners

Other restrictions imposed by the lack of implementation time were that the length of optional text elements was restricted, meaning that it was not possible to include longer sentences or paragraphs as optional elements in the document, which would have allowed us to introduce much more flexibility in adapting the document to user requirements. Also, the syntax of the option slots did not allow us to include slots within slots. So we were unable, for example, to employ a user entry slot in an option slot which would have allowed us to group all the relevant information in one slot, e.g. {one of these days/on [DATE]}. Instead we had to put an optional date slot after the option slot, e.g. {one of these days/on} [DATE], which the user would leave blank unless they selected 'on' in which case they would enter the date required. This is obviously not an ideal solution from the user's point of view since it means cluttering up the text with slots which may not be needed. It also means that the user has to read the help text and decide what to do about the date slot. For both these reasons this problem will be dealt with in the next version.

The major disadvantage with this type of application, however, is that because the user does not need to know anything about the text generated, it is likely that they will not understand a reply written in the language of the recipient. Of course a simple solution would

be to add a postscript to the letter explaining that the translated text has been generated automatically using Sharp's MDG software and requesting that a reply be written in the language of the sender if possible. An alternative possibility, if a reply is received in electronic form, is for the user to apply translation software (e.g. SID) to understand the reply. This is the approach which has been taken in the English Email Generation Assistant software which is available with the Power E/J v2.0 translation software on the Sharp Mebius 7800 Notebook PC. The user generates the English text by interacting with the system in Japanese. The text can then be translated back into Japanese if required, so that the user can review the output.

Despite these problems, we are satisfied with the performance of the system as it stands, and it has recently won an award in the New Technologies section of the "Languages for Export 97" scheme organised by the London Chamber of Commerce and Industry. The mere process of implementing it has given us many ideas and insights into the way such systems must be improved in the future. We have already carried out a detailed analysis of the existing system and have produced an initial design document for a second version of the MDG. Major improvements will include extension to further languages (including Japanese), extension of the number of document templates provided (particularly the business ones), provision for more flexible command types to permit greater flexibility between languages, access to online bilingual dictionaries, and porting the system to other hardware platforms (e.g. PC and PDA). We may also consider applying the same technology to provide the user with monolingual document templates. Interesting related research is being carried out in an EC-funded project called MABLe (Multi-Lingual Business Letter Authoring Tool, see European Union, 1997), which aims at providing templates and formulaic phrases to generate letters in the target language. However, the users of this system will need to have some knowledge of the target language, whereas the major advantage of the MDG is that the system does not require knowledge of the target language. It will be interesting to see how far it is possible to extend the concept of the MDG and provide the users with the flexibility they require, whilst at the same time maintaining 100% correct output generated automatically by the system.

## SHARP INTELLIGENT DICTIONARY (SID)

The Sharp Intelligent Dictionary (English-Japanese version) has recently been released on the Shoin word processor. It makes use of a large Sharp in-house English-Japanese dictionary developed over 15 years and which prior to this had been used as part of the Power E/J translation software. SID forms part of a suite of translation tools (which also include full text translation and phrase translation) available within the Power E/J application (Hirai, et al 1996). SID does not carry out full translation, but rather looks up words and phrases in the context of the sentence in which they are found. This context-sensitive lookup technology has been developed at Sharp Laboratories of Europe in Oxford. The major benefits for the user over existing systems are that the system provides fast, accurate lookup of words and phrases in a dictionary. The system works out the most likely translations for the words and phrases making up the sentence and then presents them to the user.

Although the system is not intended to achieve full translation – it is primarily aimed at

people who are looking at English documents they have found on the Web, who may understand some English but need help with particular words and phrases - it is interesting to note that Sharp and other vendors of MT systems, realising the difficulty of achieving accurate MT particularly with long sentences, have recently been offering translations of phrasal units as an alternative to full translation (Fukumochi, 1995). SID is perhaps better thought of as a browser or a glosser, allowing the user to scan information quickly and independently before deciding whether to download the document and possibly have it translated in full. A further use for the system might be as an aid to professional translators who could use the output of the system to assist them in producing their own translations. The *SID* user interface is shown in Figure 8 below.
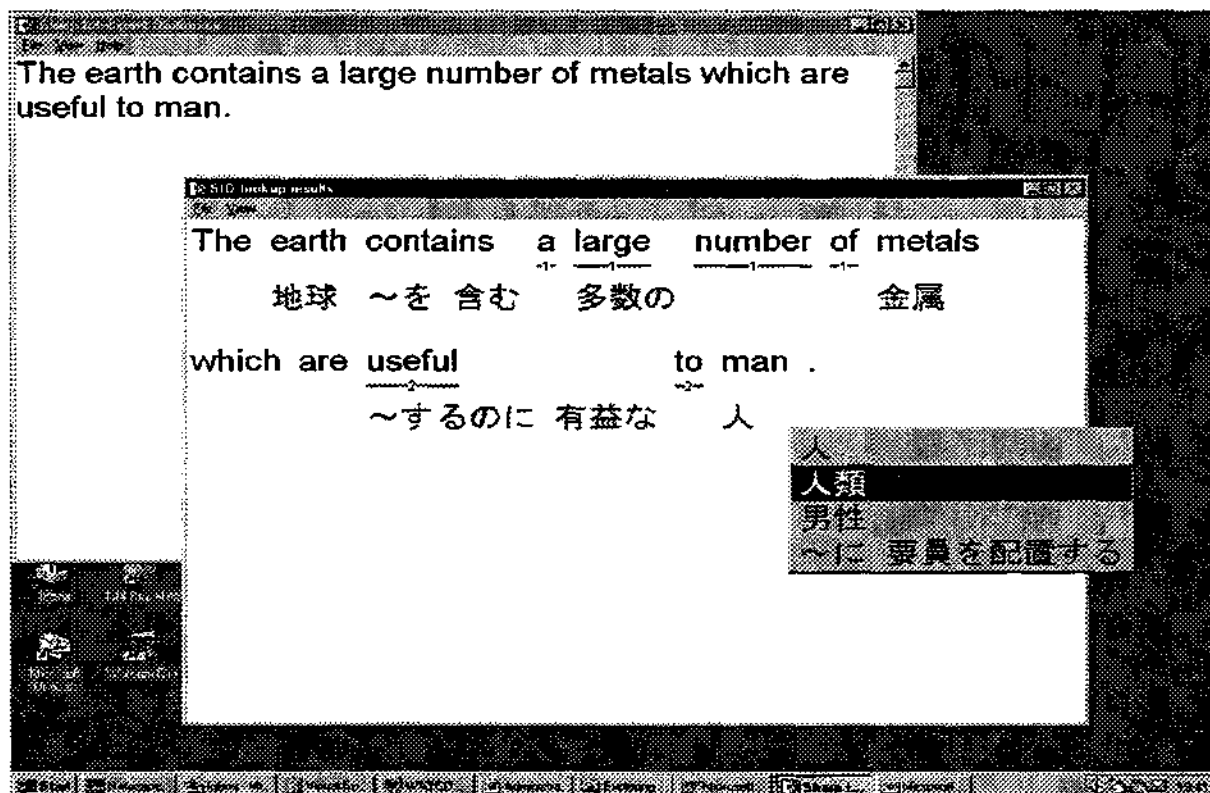


Figure 8: SID E-J showing display of alternative translations

SID works in the following way. The text is analysed by reducing words to their base forms, automatically assigning parts of speech to each word and identifying words which are part of longer, multi-word expressions. This information allows the system to retrieve the translation(s) related to a particular word or phrase in a given context from the dictionary database. The example in Figure 9 of an English-German version of SID illustrates the benefit to the user of contextual lookup over 'blind' lookup. A typical dictionary lookup system just provides a list of translation possibilities with different parts of speech: the user has to work out which one is appropriate to the particular grammatical context in which the word is being used. By contrast, SID employs knowledge of the context in which the word is being used and automatically looks up the translations for the most probable part of speech for that word. Thus, in the example below, SID correctly selects the verb 'meinen' and the adjective 'gemein' from the list of possible translations for the two uses of the word 'mean' in the context of the English sentence "I do not mean John is a mean fellow":
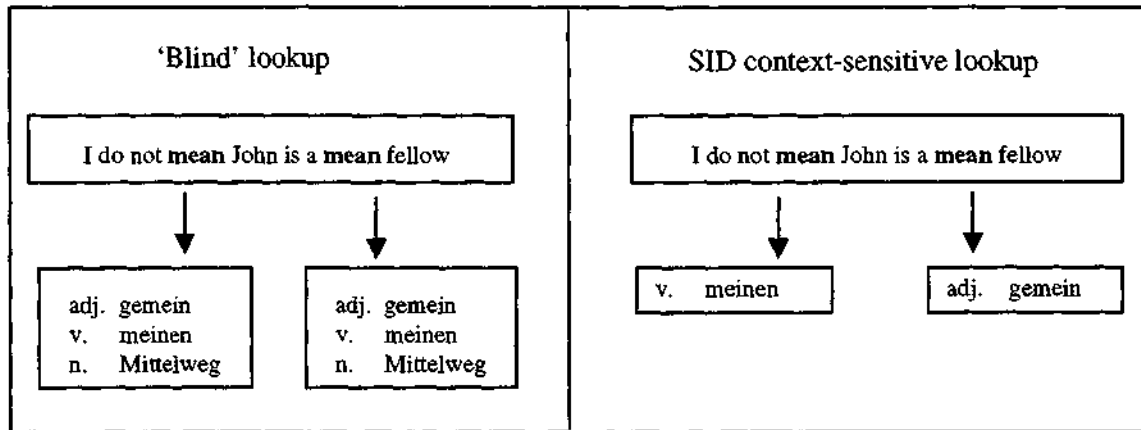
Figure 9: SID disambiguates the parts of speech of words by using contextual information

In general, the system prefers longer matches to shorter ones, so that, all other things being equal, an expression which accounts for four words in the sentence will be preferred over one that accounts for only three. Our evaluation of our current system has shown that this is generally a reliable strategy, but that in some cases the system may present an incorrect translation first. For this reason we have incorporated a function which allows the user to select an alternative from a menu accessed by clicking on the relevant word(s). This ability of the system to find multi-word expressions even if other words are inserted in between the
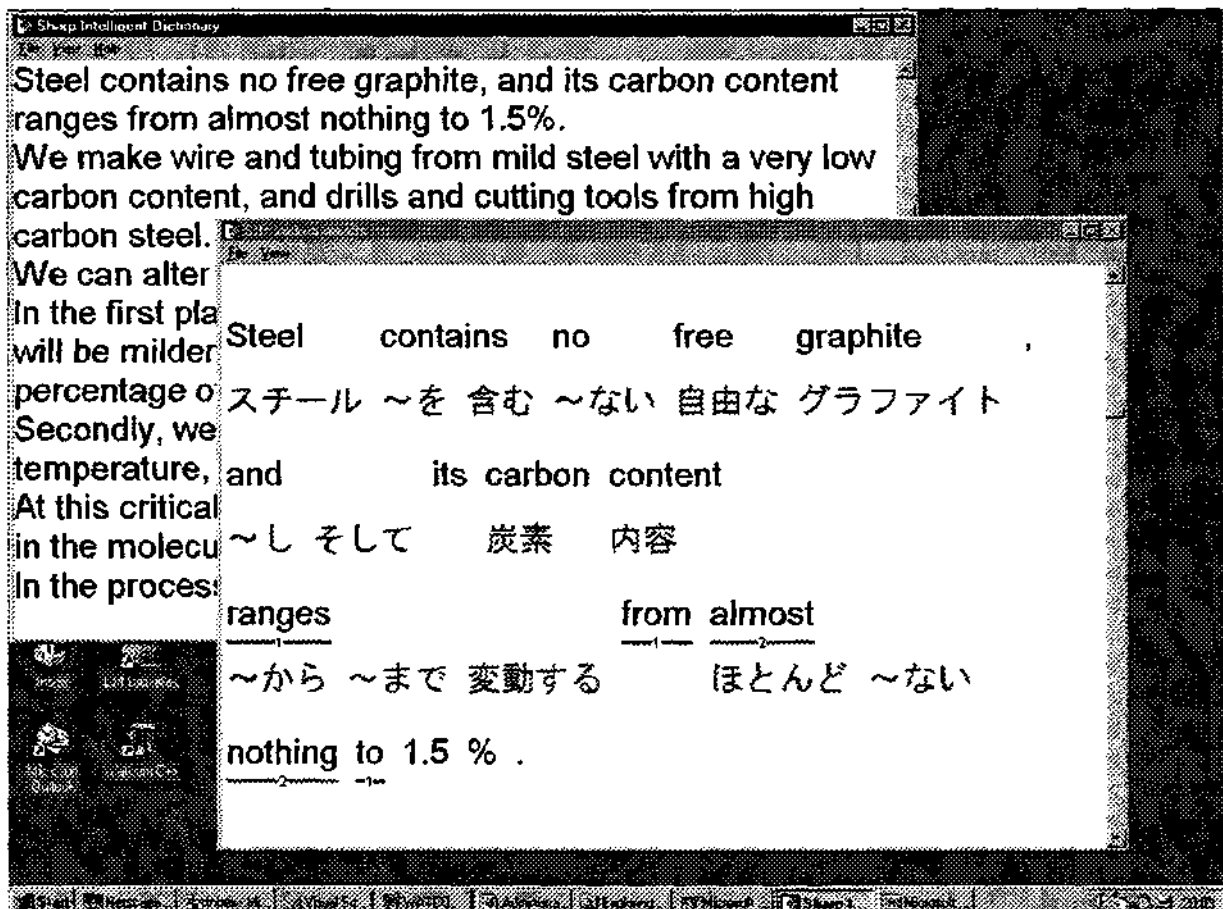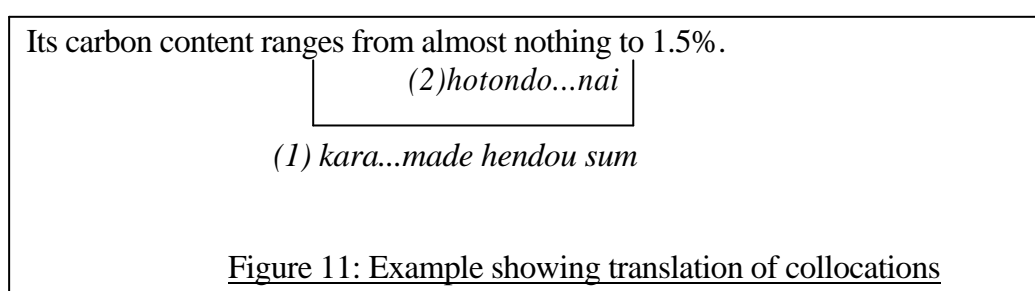


Figure 10: View of SID E-J showing display and translation of collocations

elements of the expression is one of the most useful aspects of SID. Thus, for example, SID can correctly find the multi-word expression 'make use of' in the sentence 'John is <u>making</u> considerable <u>use of</u> the Internet', despite the insertion of the word 'considerable'. This is an improvement on the performance of some other dictionary accessing software (e.g. the thesaurus available in MS Word 7.0), which is unable to find this expression when other words are placed in between, although the expression is known to the system. We are planning to introduce this feature in future versions of the product.

A particular feature of SID is the identification of collocated words by means of underlining in colour with numbers grouping the components of a collocate, as illustrated in the following example taken from Figure 10 above:

Its carbon content ranges from almost nothing to 1.5%.
*(2)hotondo...nai*

*(1) kara...made hendou sum*

Figure 11: Example showing translation of collocations

We have also improved our part-of-speech tagger to enable it to distinguish between transitive and intransitive verbs. This allows SID to provide the user with more accurate translations, as illustrated in the following example:

| **Previous version of SID** | | | **Latest version of SID** | | |
|---|---|---|---|---|---|
| He | runs | <u>for his life</u>. | He | runs | <u>for his life</u>. |
| | *hashiru* | *inochigakede* | | *hashiru* | *inochigakede* |
| He | runs | a program. | He | runs | a program. |
| | *hashiru* | *puroguramu* | | *hashiraseru* | *puroguramu* |

Figure 12: Improved treatment of transitive/intransitive verbs

We have plans to improve all aspects of the functionality of the system including the graphical user interface, the analysis component and the dictionary coverage. In the future we would like to extend the tagging software in the direction of the identification of the senses of individual words and phrases as well as their grammatical parts of speech. We would like the system to be able to identify the particular sense of the word 'bank' (financial institution, river edge, tilt aircraft to one side) as well as the part of speech (noun or verb). This would enable us to increase further the accuracy of the dictionary lookup mechanism. Some initial research has been carried at Cambridge Language Services (CLS) and at SLE in this area (see paper by Harley and Glennon at the ACL Workshop on Semantic Tagging 1997), principally in the context of our collaboration on the Integrated Language Database project (1993-7) which has been supported by partial funding from the DTI and EPSRC. Early results are encouraging

and we hope to take them account in improving our current system. Finally, although SID is currently available for English-Japanese only, we have plans to extend the range of language pairs in the near future. Initial pairs will be from English into the target language but we are planning to carry out experiments on applying the SID technology to other source languages.

## CONCLUSION

Sharp is committed to providing a range of personal translation tools ranging from document generation to lookup of words and phrases in a variety of dictionaries in order to help users approach foreign languages with greater confidence. The degree of linguistic competence expected of users ranges from zero (in the case of the MDG) to some basic knowledge of English (in the case of SID). People with more advanced knowledge are also catered for with the MDG in that even if one has a very good knowledge of foreign languages it is very difficult always to know the correct stylistic way of expressing oneself, particularly for example in the case of delicate situations (e.g. a letter of sympathy about a bereavement). The MDG provides templates which are guaranteed correct both linguistically and stylistically. Similarly, SID may be useful even to those with an advanced knowledge of English, because the system allows the selection of different user levels, thus ensuring that advanced users are not shown translations for relatively simple words. The system also has a variety of technical dictionaries which have been developed in the course of these R&D activities.

By providing personal translation products such as those described in this paper which allow people to browse documents on the Internet in a language they may not know or generate standard letters without knowing a word of the target language, Sharp is carrying out its mission to improve global information access and communication for individuals living in the Information Age.

## ACKNOWLEDGEMENTS

# REFERENCES

1.  Davies, S., et al, 1989, <u>Bilingual Handbook of Business Correspondence and Communication</u> (Prentice Hall International)

2.  Femey, D., et al, 1990, <u>The Multilingual Business Handbook</u> (Macmillan)

3.  European Union Telematics Applications Programme (Language Engineering Sector), 1997, <u>Language Engineering in Europe: Progress and Prospects</u> (LINGLINK, Anite Systems)

4.  Fukumochi, Y., 1995, <u>A Way of Using a Small MT System in Industry</u> (Proceedings of MT Summit V, Luxemburg)

5.  Harley, A., and Glennon, D., 1997, <u>Sense Tagging in Action</u> (ACL-97 Workshop on Semantic Tagging)

6.  Harvard, J., et al, 1989, <u>Bilingual Guide to Business and Professional Correspondence</u> (Pergamon Press)

7.  T. Hirai, et al, April 1996, <u>English-Japanese Translation Support Software Adapted for the Internet,</u> Sharp Technical Journal Vol 64, April 1996

8.  Jones, D. B., and Sager, J. C., 1989, <u>Automatic Translation Perspectives for Letter-Writing Systems</u> (UMIST. Centre for Computational Linguistics Report No: TE89-1)

9.  Saito, H., and Tomita, M., 1985, <u>On Automatic Composition of Stereotypic Documents in Foreign Languages</u> (Carnegie-Mellon University, Department of Computer Science Technical Report CMU-CS-86-107)