

# Sharing User Dictionaries Across Multiple Systems with UTX-S

AAMT Sharing/Standardization Working Group,

Francis Bond,<sup>1</sup> Seiji Okura,<sup>2</sup> Yuji Yamamoto,<sup>3</sup> Toshiki Murata,<sup>4</sup>  
Kiyotaka Uchimoto,<sup>1</sup> Michael Kato,<sup>5</sup> Miwako Shimazu,<sup>6</sup> Tsugiyoshi Suzuki<sup>7</sup>

<sup>1</sup> National Institute of Information and Communications Technology,  
<sup>2</sup> Fujitsu Laboratories Ltd., <sup>3</sup> CosmosHouse, <sup>4</sup> Oki Electric Industry Co., Ltd.,  
<sup>5</sup> Learning Consultant, <sup>6</sup> Toshiba Solutions Corporation, <sup>7</sup> Cross Language Inc.

## ABSTRACT

Careful tuning of user-created dictionaries can be very advantageous when using a machine translation system for computer aided translation. However, there is no widely used standard for user dictionaries in the Japanese/English machine translation market. To address this issue, AAMT (the Asia-Pacific Association for Machine Translation: <http://www.aamt.info/>) has established a specification of sharable dictionaries (UTX-S: Universal Terminology eXchange — Simple), which can be used across different machine translation systems, thus increasing the interoperability of language resources. UTX-S is simpler than existing specifications such as UPF and TBX. It was explicitly designed to make it easy to (a) create new user dictionaries and (b) share existing user dictionaries. This facilitates rapid user dictionary production and avoids vendor tie in. In this study we describe the UTX-Simple (UTX-S) format, and show that it can be converted to the user dictionary formats for five commercial English-Japanese MT systems. We then present a case study where we (a) convert an on-line glossary to UTX-S, and (b) produce user dictionaries for five different systems, and then exchange them. The results show that the simplified format of UTX-S can be used to rapidly build dictionaries. Further, we confirm that customized user dictionaries are effective across systems, although with a slight loss in quality: on average, user dictionaries improved the translations for 44.8% of translations with the systems they were built for and 37.3% of translations for different systems. In ongoing work, AAMT is using UTX-S as the format in building up a user community for producing, sharing, and accumulating user dictionaries in a sustainable way.

## Author Keywords

Interoperability of Language Resources, Machine Translation, User Dictionaries

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

*IWIC'09*, February 20(21, 2009), Palo Alto, California, USA.

Copyright 2009 ACM 978-1-60558-198-9/09/02...\$5.00.

## ACM Classification Keywords

I.2.7 [Artificial Intelligence] Natural Language Processing—Machine translation, I.2.m [Artificial Intelligence] Miscellaneous

## INTRODUCTION

User and specialist domain dictionary compilation is important for machine translation, because adding terms into a user dictionary decreases failures in syntax analysis [5] as well as improving the accuracy of translation. Yoshimura *et al.* [17] using 41 hierarchically organized domains (which they call Subject Areas) marked on 9,000 words, were able to improve the translations of 12% of the badly translated nouns ( $\frac{46}{383}$ ). Similarly, Lange and Yang [10] used 77 Domains (TERMinology CATegories) and 30 Topical Glossaries. In a first pass through they chose two appropriate domains (the domains with the greatest number of domain-tagged words). With these domains set, there were changes in 0–40% of translations, and the majority of the changes were improvements.

Another measure of how valuable user dictionaries are is the lengths people are prepared to go to make them work with existing systems. Inaba *et al.* [6] develop a "Specialized dictionary" that is applied to an MT system accessed only through a server. Words in the specialized dictionary are replaced by special words (e.g. "REF1") and their translation is then replaced by the translation of the specialized word from the dictionary. This allows the use of external user dictionaries, even when the MT system is a black box. Itagaki and Aikawa [8] use a similar technique to replace words in an SMT system — in their case translating the source word in various contexts, finding how the system translates it and storing that in a preprocessing step. Then in actual translation, they replace the system's translation of the source word with the translation in the user dictionary. In both these cases considerable effort has been expended in retrofitting user dictionaries to existing MT systems.

However, user dictionary compilation is a very time-consuming process for an individual user, and its effect may not immediately be clear for relatively small translation projects. In addition, the formats of individuals' dictionaries are varied

if they do not use the same machine translation software. Therefore, sharing these dictionaries can be difficult.

There are already several well designed formats for user dictionaries/terminology: TermBase eXchange: **TBX**<sup>1</sup>; the Open Lexicon Interchange Format: **OLIF**<sup>2</sup> [11] and the Universal PlatForm: **UPF**<sup>3</sup> [9]. However, to the best of our knowledge, there are not many dictionaries available formatted as either TBX, OLIF or UPF, although all are supported by some commercial machine translation systems.

More recently the general Lexical Markup Framework (**LMF**) has also been proposed [4] as a framework for representing dictionaries. However, LMF is a framework, it has to be specialized further in actual use, and shares the complexity of the above systems.

All these specifications are intended to be used by professional terminologists. Their complex nature gives them the ability to represent the various kinds of information needed in a machine translation system. However, this makes them hard to use for average MT users, and none of them are widely used in Japan. These concerns have been noted before, leading to a simplified format, **TBX basic** [12], which is considerably streamlined, but it still requires a great deal of information to be entered beyond simple translation equivalents. The majority of MT users in Japan still use simple text formats or spreadsheet lists. In order to address these problems, we have established a new standard called UTX-Simple (**UTX-S**). It is not meant to replace any of the existing standards for heavy machine translation users, but is instead meant to supplement them as a lightweight variant. As far as possible, we are making the standard compatible with the more detailed standards.

In this paper, we introduce UTX-S. We then use UTX-S in two experiments: one on reusing domain dictionaries and one on sharing user dictionaries. Finally we discuss future work.

## UTX-SIMPLE

The features of UTX-S are as follows:

1. **Dictionary for the user - simple and easy to use:** A complicated specification merely increases users' burdens, and it will be forgotten eventually. A specification must be simple and practical to reflect and include actual MT users' needs, viewpoints, and scenarios.
2. **Entry as a technical term:** UTX clearly defines the domain of a dictionary, and adheres to the principle of one word, one meaning. Because the format does not allow one to enter information to choose between translations, a given term is restricted to a single translation. Therefore, the domain should be restricted enough that an entry should be a unique term within an applicable domain.

<sup>1</sup><http://www.lisa.org/Term-Base-eXchange.32.0.html> [16]

<sup>2</sup><http://www.olif.net/index.htm>

<sup>3</sup><http://www.aamt.info/japanese/upf.htm>

3. **Improvement in translation accuracy:** When UTX dictionaries become widespread, sharing dictionaries can be drastically simplified. Users can compile user dictionaries more efficiently by exchanging existing data, and UTX will eventually contribute to improvement in translation accuracy.
4. **Multilingual and monolingual dictionary:** UTX is designed for dictionaries not only between two languages but also among multiple languages. In addition, UTX also includes a specification for a monolingual dictionary, which can be used for proofreading tools for terminological standardization, etc.
5. **Promotion of localization of software:** Especially in open source localization projects, translation is carried out individually, and terminological standardization can be difficult. By exchanging and sharing translation resource systematically through the use of UTX, more and more common dictionaries become available, and they increase the efficiency of translation exponentially.

As UTX-Simple does not retain detailed information about its creator, the timestamp for each entry, etc., it is not suitable for a compilation of a permanent, versatile dictionary. However, it is practical and easy to create, use, and share. It may be used by a variety of users regardless of their background knowledge in advanced linguistics.

We are currently establishing guidelines for notation, i.e. how to describe each entry, for each target language. As far as possible we will try to stay compatible with existing guidelines, such as those for TBX, OLIF and JMDict [2].

The converter which maps between UTX-S and MT system formats also acts as a validator. For example, if one or more mandatory properties are missing, a warning is issued according to the corresponding option of the tool. To cope with human errors, the tool will skip illformed lines during conversion. However, by making more detailed guidelines, and also a format defined well enough to be validated, we hope to make it easier for people to make their existing lexical resources useful for other people.

## UTX-S Specifications

A UTX-S dictionary file is a plain text file, in UTF-8 encoding. It consists of a header (the first two lines) and the body (all subsequent lines). We give an example in Table 1. The full specifications for UTX-S are online at <http://www.aamt.info/english/utx/>.

The first header line consists of the following fields, separated by semicolons:

- An initial #
- The string UTX-S and the version number (currently 0.92)
- The source language (expressed as ISO 639-1 language and ISO 3166-1 alpha-2 country codes)
- The target language (same as source)
- The date produced

- Misc (normally creator and copyright info)

The second header line consists column headers, separated by tabs, the first three columns are mandatory.

- An initial #
- **src**: the source language word
- **tgt**: the target language word
- **src:pos**: the source part of speech  
{noun |verb |adjective |adverb |properNoun |sentence}
- zero or more user defined columns, e.g.:
  - tgt:pos
  - src:plural
  - comment

Comment lines (defined as any line starting with a hash #) can be freely interspersed with the dictionary entries. This means that it is easy to delete entries (just comment them out).

Note that the 4th and subsequent fields are all optional: in the experiments described below we only used the first three mandatory fields. The three mandatory fields are the same as for TBX-basic, and almost the same as those described in use by Dillinger [3] for user dictionaries built for Logos (source term, source gender, source part of speech and target term). Gender is not grammaticalized for either Japanese or English, the two languages we are most interested in, so we do not make it a mandatory field. It could be added, for example as `src:gender`.

### UTX Entry Guidelines

In this system we give the guidelines for writing individual entries.

- Add only one translation for each word-pos pair.
- Avoid words already in the system dictionary.

#### Part of Speech

- Choose one from the following:
  - noun
  - verb
  - adjective
  - adverb
  - properNoun
  - sentence
- If the POS is unknown then leave it blank: "".

The “sentence” part of speech is for complete utterances that should be translated as is, with no further processing. It is not meant to replace a full translation memory, but rather is made available for things like greetings and proverbs, which are rarely compositional in translation, and often hard to parse. This part-of-speech is supported by most of the commercial MT systems we tested.

#### English Entry

- Use the singular form by default
- Use lower case letters
- Do not add articles (*a*, *an*, *the*) or infinitival *to*.
- Do not use full-width letters

#### Japanese Entry

- Add `する -suru` to nominal verbs when translated as a verb. E.g.,
 

翻訳	translation	noun
翻訳する	translate	verb
- Add `な -na` to nominal adjectives
- noun-no combinations can be treated as adjectives  
e.g.: orthographic 綴りの adjective
- Use either a middle dot or a space to separate foreign compounds. E.g.,  
ピタビ・アルゴリズム Viterbi algorithm noun

### EXPERIMENTS

UTX is designed to make it simple to (a) create new user dictionaries and (b) share existing user dictionaries.

In order to test the applicability of the UTX-S guidelines, we performed two experiments. In the first we converted an existing specialist dictionary to UTX-S, and in the second we compared user dictionaries across different systems.

#### Test Data

For both experiments we tested by translating a small sample document. We used the OLIF Guidelines for Formulating Canonical Forms, a 147 sentence English document from the on-line documentation for OLIF,<sup>4</sup> and translated it into Japanese. The text was extracted and sentence aligned manually — the input to the systems was clean, sentence separated text. We give an example, with a manually constructed reference translation and four machine translation outputs in Figure 1. We consider this to be text in the domain of natural language processing.

Testing was done by (a) regression testing, as this has been shown to be consistent and fast [14] and with (b) BLEU (with up to 4 n-grams) [13]. The baseline was the input text translated by the MT system with no user dictionary.

For the regression test, the results of using different dictionaries were compared to the original translation using only the system dictionaries. Translations were presented to an evaluator in random order, so that the evaluator did not know which system was which. Identical translations were scored automatically. Changed translations were evaluated as either equivalent quality (0), improved (+1) or degraded (-1). The total change was calculated by summing the scores.

All judgments were made by only one evaluator, with a different evaluator for each system (for a total of five different people). Because of the known reliability of regression

<sup>4</sup><http://www.olif.net/documentation.htm>

#UTX-S 0.92; en-US/ja-JP; 2008-03-15T10:00:00Z+09:00; copyright: AAMT, license: CC-by 3.0									
#src	tgt	src:pos	src:plural	src:3sp	src:past	src:pastp	src:presp	src:comp	src:super
new	新規の	adjective						newer	newest
fast	高速な	adjective						faster	fastest
# prosody should be uncountable									
prosody	韻律	noun	prosodies						
save	保存する	verb		saves	saved	saved	saving		
good evening	今晩は	sentence							

Table 1. Example of UTX-S

• *OLIF Guidelines for Formulating Canonical Forms*

- (a) OLIF・正規化形式への定型化の指針 (reference)  
“OLIF Guidelines for Formulating Canonical Forms”
- (b) 教会法に基づく形式を定式化するためのOLIFガイドライン (MT)  
“OLIF Guidelines for formularizing forms based on Canon Law”
- (c) 教会法に基づくフォームを定式化するためのOLIFガイドライン (MT+lingdic)  
“OLIF Guidelines for formularizing forms based on Canon Law”
- (d) 正規化形式を形式化するためのOLIFガイドライン (MT+user)  
“OLIF Guidelines for formulating canonical forms”
- (e) 基準形を定型化するためのOLIFガイドライン (MT+other)  
“OLIF Guidelines for formulating regular forms”

Figure 1. Example Translation with different User Dictionaries

testing, and the small scale of the project, we did not use multiple evaluators on the same data.

BLEU testing was done with a single reference translation, prepared by a professional translator with expert domain knowledge. In BLEU the system translation is compared to the reference (human translation) by measuring the number of n-grams which overlap between the two translations. We used `multibleu.perl` comparing up to 4-grams.

For example, consider sentence (b) *OLIF Guidelines for formularizing forms based on Canon Law* compared to the human reference translation (a) *OLIF Guidelines for Formulating Canonical Forms*. There are three 1-grams that overlap (*OLIF*, *Guidelines*, *for*), two 2-grams (*OLIF Guidelines*, *Guidelines for*) and one 3-gram (*OLIF Guidelines for*).

The test set is small (147 sentences), but the automatic evaluation showed the same trends as the human evaluation.

*Systems*

Five commercial MT systems were tested:

- LogoVista PRO 2008 Super Pack  
([www.logovista.co.jp/product/honyaku\\_pro2008/pro2008\\_st.html](http://www.logovista.co.jp/product/honyaku_pro2008/pro2008_st.html))
- Translation Software ATLAS  
([www.fujitsu.com/global/services/software/translation/atlas/lineup/](http://www.fujitsu.com/global/services/software/translation/atlas/lineup/))
- Collaborative Translation Environment: Yakushite.Net  
([yakushite.net/](http://yakushite.net/): [15])

- PC-Transter 2008 Professional (Cross Language Inc.)  
([www.crosslanguage.co.jp/products/studio2008/](http://www.crosslanguage.co.jp/products/studio2008/))

- The HON-YAKU 2008 Premium  
([pf.toshiba-sol.co.jp/prod/hon\\_yaku/premium/index\\_j.htm](http://pf.toshiba-sol.co.jp/prod/hon_yaku/premium/index_j.htm))

Results are given for the five systems anonymized as **A**, **B**, **C**, **D** and **E**.

**Converting an Existing Lexicon**

In the first experiment, we wanted to test how useful an off the shelf domain dictionary was as a user dictionary. We made a user dictionary based on **lingdic**. **lingdic** is a small open source dictionary of Japanese-to-English computational linguistic terms.<sup>5</sup> It is maintained by a single volunteer, enhanced with contributions from the community. It is formatted based on the format used by EDICT, a large collaborative Japanese-English dictionary [1].

The version we used in our experiment has 3,527 Japanese head words and 4,123 Japanese-English pairs. We wrote an `edict2utx-s` converter (including a mapping from the edict POS set to the UTX-S set), and enhanced the dictionary in the following ways:

- Added more part-of-speech information, in particular, tags for all verbal-nouns and nominal-adjectives.

<sup>5</sup>[www2.nict.go.jp/x/x161/en/member/bond/data/lingdic](http://www2.nict.go.jp/x/x161/en/member/bond/data/lingdic)

- Added a tag (**xmt**) to mark pairs which should not be used in a machine translation dictionary.
- Produced a reverser — a program to recast the dictionary as an English-to-Japanese dictionary.

The English - Japanese dictionary is made by reversing each entry and sorting by the following criteria, cascaded.

1. Prefer similar translations, more specifically

- prefer acronym↔acronym, long form↔long form
- |                                      |           |   |
|--------------------------------------|-----------|---|
| HPSG                                 | HP S G    | ↑ |
| HPSG                                 | 主辞駆動句構造文法 | ↓ |
| head-driven phrase structure grammar | 主辞駆動句構造文法 | ↑ |
| head-driven phrase structure grammar | HP S G    | ↓ |

2. Prefer the translation with the highest monolingual frequency  
(frequencies were obtained using the Yahoo Japan Developer API).

bottom-up	ボトムアップ	1,020,000	↑
bottom-up	ボトムアップ法	468,000	↓

3. Prefer the shorter translation

4. Sort in unicode order (to make the sort deterministic)

Note that the target language frequency is rarely the same, so 3. and 4. are almost unused.

Some example entries for the converted, reversed **lingdic** are given in the appendix, in Table 3. In our first experiment, this dictionary was converted from UTX-S to the native formats of the five MT systems, and used as a user dictionary. The results for this configuration are given as +**lingdic** in Table 2.

### Creating and Exchanging User Dictionaries

In our second experiment, five members of the AAMT working group (one for each system) translated the document once, and then created a user dictionary designed to fix some of the problems in the MT results (to be combined with **lingdic**). This is not a blind test, in that the users see the machine translation output. However, this is how users typically create their own user dictionaries. Each person spent a small amount of time preparing the user dictionary (generally a couple of hours). The number of words added ranged from 17–156, as shown in Table 2.

There was a large amount of variance in which words were added. Only one word was added to all five user dictionaries: *compound* which was entered as ⟨compound, 複合語, noun⟩ *fukugougo* “compound expression” in four user dictionaries and ⟨compound, 複合名詞, -⟩ *fukugou-meishi* “compound noun” in one. There were also many places where different lexical choices were made: for example *string* was translated as ⟨string, 文字列, noun⟩ *mojiretsu* “character array” by three users and ⟨string, ストリング, noun⟩ *sutoringu* “string” by two users. Either is acceptable, although use should be consistent within a document.

One of the testers also added some entries that we would expect to be specific to this text, for example ⟨write- as a stem for write, writeの語幹のwrite-, noun⟩.

The results for adding this customized user dictionary to **lingdic** and translating are given as +user in Table 2.

For the final experiment, we tested the hypothesis that a user dictionary built for one system will also be useful in a different system. If this were not the case, then there would be no motivation to make an interoperable user dictionary format. To test this we simply exchanged user dictionaries: System A uses the dictionary created for system E, B uses the one for A, C uses the one for B and so on. The results for using this combined user dictionary are given as +other in Table 2.

Finally, the translator who made the reference translation made a user dictionary by merging all five dictionaries, and adding a few more entries. We used this for one final experiment to see just how good we could get using only user dictionaries.

## RESULTS AND DISCUSSION

The results of the three experiments are shown in Table 2.

The result of the first experiment, adding a reversed existing specialist dictionary (+**lingdic**), was negative. The dictionary had a negative effect on most of the systems, and only a very small positive effect on one. On average, 6.5% of sentences had their translation degraded, and the BLEU score went down by 0.6. There were two main reasons for the degradation: the first was that most of the systems (all but E) had an over-strong preference for words in the user dictionary — they were preferred even over existing multiword expressions. An example of an error caused by this was “upper case” going from 大文字 *oomoji* “capital letters” to 上の格 *ue-no-kaku* “upper (grammatical) case”. **lingdic** had an entry for 格 *kaku* (grammatical) “case” (which was also used correctly in the document), but it caused a degradation for this term. This was corrected in the hand-built user dictionaries by adding an entry for ⟨upper case, 大文字, noun⟩ and even for the larger noun phrase ⟨upper and lower case, 大文字と小文字, noun⟩.

The second reason was that the best word was often not chosen when the dictionary was reversed. For example *word* was translated as 用語 *yougo* “term, word” rather than 単語 *tango* “word”. 用語 is more common, but is a less good translation for *word* in the linguistic domain. Instead of monolingual frequencies, we really need the frequencies conditioned on the source word, although these can only be found from bilingual corpora. An alternative would be to try and collect monolingual target language documents from the domain in question (computational linguistics/NLP) and use them to calculate the frequency, rather than the more general Web as corpus. Deleting only a few entries made an enormous difference, although which words were problematic varied from system to system.

System	Dic Size	Percent Change			BLEU Score			
		+ling	+user	+other	System	+ling	+user	+other
A	156	-1.4	66.0	38.1	13.8	13.2	21.4	18.2
B	59	-19.7	56.0	20.4	15.4	15.8	18.6	21.1
C	27	-5.4	27.2	36.7	15.5	14.7	17.6	17.0
D	24	-8.2	45.6	32.7	17.2	15.3	20.3	17.8
E	17	2.0	46.9	58.5	12.2	11.7	16.5	16.4
Ave	56.6	-6.5	44.8	37.3	14.8	14.2	18.9	18.1

+ling: results with **lingdic-EJ**

+user: results with a user dictionary built for that system

+other: results with an exchanged user dictionary: A uses E, B uses A, C uses B and so on.

**Table 2. Translation Results**

Dillinger [3] also found that converting domain dictionaries did not always produce a useful dictionary for MT (in particular many entries in a specialist dictionary were never actually encountered in text). However, people found the dictionary useful when building their own user dictionaries, and UTX-S makes it was easy to add new entries and delete existing ones by commenting them out. One of the goals of the AAMT working group is to try not only to convert existing dictionaries to UTX format, but also to validate them in use and get feedback. In this case, by adding the **xmt** tag (don't use for MT) to the lingdic source, we were able to make the human use dictionary more useful for machine translation.

The results of creating a user dictionary tuned to the text showed, as expected, that it improved the translation quality (+user). Translations improved for 44.8% of sentences according to the regression testing and the BLEU score went up by 4.1. The BLEU score almost certainly under counts the improvement. Because there is only one reference translation, legitimate variations such as those for *string*, where the native Japanese 文字列 or the loan word ストリング are pretty much interchangeable, are penalized.

We were not able to quantify the ease-of-use of the format, but qualitatively all the testers found it easy to use. In particular, people liked being able to edit the dictionaries in the editor/spreadsheet of their choice. It was definitely easier to share than proprietary dictionary formats.

In the next experiment, we found that user dictionaries, even if developed for other systems, improved the translation quality (+other). As far as we know this has been assumed before, but never directly shown. Translations improved for 37.3% of sentences and the BLEU score went up by 3.3 compared to using no dictionary.

Translations for all four configurations, plus the reference translation, are shown in Figure 1, with glosses. The best automatic translation is (d) with the system-specific user dictionary, but (e) with the dictionary built for another system is also perfectly understandable.

Finally, we calculated the BLEU score for the translation with the merged, corrected dictionary (using system D). It was 44.52, an improvement of 27.3 points. This shows just

how useful user dictionaries can be. However, this dictionary has 146 entries, almost one per translated sentence, so it is beyond what would normally be feasible.

### Dictionary Conversion Tools

We have created a perl script that converts from UTX-S to the user dictionary formats of the five systems tested here. The conversion back from user dictionary to UTX-S would be lossy: UTX-S is less informative than the full user dictionary formats. The script is available from the UTX page (<http://www.aamt.info/english/utx/>).

### FUTURE WORK

Having confirmed the ease-of-use of UTX-S and the efficacy of cross-system user dictionaries, we now to turn to investigating methods of encouraging people to make and share their dictionaries.

Our basic plan is to produce or encourage the production of domain specific dictionaries, such as education, sport, IT, medicine etc. These dictionaries would be tested with MT systems, removing (or marking) entries that may be useful for humans, but which degrade translation quality. We will not work on a general lexicon, as all MT systems come with their own system dictionaries.

In open source localization projects, translation is carried out individually, and dictionaries are not shared as they should be. Dictionaries are scattered across various providers, and their licenses and formats are also varied. If these scattered language resources are centralized, the localization between different languages can be significantly accelerated.

In order to spread UTX dictionaries in which anyone can participate in creation, and in order to realize open dictionaries for everyone, a shared dictionary community should be established. We are considering two types of dictionary communities for producing, sharing, and accumulating dictionaries. They can be distributed using a common framework.

- The official dictionary community (managed by AAMT or its delegate) will offer validated dictionaries with guaranteed quality for a fee.
- The open dictionary community offers free dictionaries

with open source license and promotes mutual exchange. AAMT or its delegate provides hosting service only, but no management or guarantee.

In order to meet various needs, official dictionaries and free dictionaries should be distinguished. Free dictionaries can be used for no fee, although the correctness of the contents is not guaranteed.

Collaborations with other parties has already begun for UTX, including cooperation with Oki's community-oriented machine translation site Yakushite-net [15] in which users can add their own technical terms. We have also fed back improvements to **lingdic** and the JMDict project and are cooperating with other projects which connect many dictionaries and systems (such as the Language Grid [7]). While providing UTX to other parties and verifying its effectiveness, we are also interested in collaborating in terms of tool development.

## CONCLUSIONS

We have presented a case study where we use the UTX-S format to (a) convert an on-line glossary to UTX-S, and (b) produce user dictionaries for five different systems, and then exchange them. The results showed that the simplified format of UTX-S can be used to rapidly build dictionaries. Further, we confirmed that customized user dictionaries are effective across systems, although with a slight loss in quality: on average, user dictionaries improved the translations for 44.8% of translations with the systems they were built for and 37.3% of translations for different systems.

## REFERENCES

1. BREEN, J. W. Building an electronic Japanese-English dictionary. Japanese Studies Association of Australia Conference ([http://www.csse.monash.edu.au/~jwb/jsaa\\_paper/hpaper.html](http://www.csse.monash.edu.au/~jwb/jsaa_paper/hpaper.html)), 1995.
2. BREEN, J. W. JMDict: a Japanese-multilingual dictionary. In *Coling 2004 Workshop on Multilingual Linguistic Resources* (Geneva, 2004), pp. 71–78.
3. DILLINGER, M. Dictionary development workflow for MT: Design and management. In *MT Summit VIII* (Santiago de Compostela, 2001), pp. 83–88.
4. FRANCOPOULO, G., GEORGE, M., CALZOLARI, N., MONACHINI, M., BEL, N., PET, M., AND SORIA, C. Lexical markup framework (LMF). In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC 2006)* (Genoa, Italy, May 2006).
5. FUJI, M. Experiments to evaluate the reading comprehension of English-to-Japanese machine-translated texts. In *2nd Annual Meeting of the Association for Natural Language Processing* (Tokyo, 1996), pp. 21–24.
6. INABA, R., MURAKAMI, Y., NADAMOTO, A., AND ISHIDA, T. Multilingual communication support using the language grid. In *International Workshop on Intercultural Collaboration (IWIC2007)* (2007), vol. 4568 of *Lecture Notes in Computer Science, Vol. 4568*, Springer, pp. 118–132.
7. ISHIDA, T. Language grid: An infrastructure for intercultural collaboration. In *IEEE/IPSJ Symposium on Applications and the Internet (SAINT-06)* (2006), pp. 96–100. (keynote address).
8. ITAGAKI, M., AND AIKAWA, T. Post-MT term swapper: Supplementing a statistical machine translation system with a user dictionary. In *Proceedings of the Sixth International Language Resources and Evaluation (LREC'08)* (Marrakech, Morocco, 2008), E. L. R. A. (ELRA), Ed.
9. KAMEI, S., ITOH, E., FUJII, M., HIRAI, T., SAITOH, Y., TAKAHASHI, M., HIYAMA, T., AND MURAKI, K. Sharable formats and their supporting environments for exchanging user dictionaries among different MT systems as a part of AAMT activities. In *MT Summit VI* (1997).
10. LANGE, E. D., AND YANG, J. Automatic domain recognition for machine translation. In *Machine Translation Summit VII* (Singapore, 1999), pp. 641–645.
11. LIESKE, C., MCCORMICK, S., AND THURMAIR, G. The open lexicon interchange format (OLIF) comes of age. In *MT Summit VIII* (Santiago de Compostela, 2001), pp. 211–216.
12. MELBY, A. K. TBX-Basic: Translation-oriented terminology made simple. *Revista Tradumàtica (in Press)* (2009). (ISSN 1578-7559).
13. PAPINENI, K., ROUKOS, S., WARD, T., AND ZHU, W. J. BLEU: a method for automatic evaluation of machine translation. In *40th Annual Meeting of the Association for Computational Linguistics: ACL-2002* (2002), pp. 311–318.
14. PINKHAM, J., AND SMETS, M. Machine translation without a bilingual dictionary. In *Ninth International Conference on Theoretical and Methodological Issues in Machine Translation: TMI-2002* (Keihanna, Japan, 2002), pp. 146–156.
15. SUKEHIRO, T., KITAMURA, M., AND MURATA, T. Collaborative translation environment 'Yakushite.Net'. In *Proceedings of the Sixth Natural Language Processing Pacific Rim Symposium: NLPRS-2001* (Tokyo, 2001), pp. 769–770.
16. THURMAIR, G. Exchange formats: TBX, OLIF and beyond. *LDV-Forum 21*, 1 (2006), 45–56.
17. YOSHIMURA, Y., KINOSHITA, S., AND SHIMAZU, M. Processing of proper nouns and use of estimated subject area for web page translation. In *Seventh International Conference on Theoretical and Methodological Issues in Machine Translation: TMI-97* (Santa Fe, 1997), pp. 10–18.

## APPENDIX

### EXAMPLE DICTIONARIES

Some entries from **lingdic** with the word *lexicon* or *dictionary* are shown in Table 3. The user dictionary created for system D is shown in Table 4.

#UTX-S 0.92; en-US/ja-JP; 2008-05-21; copyright: Francis Bond (2008); license: CC-by 3.0		
#src	tgt	src:pos
Japan Electronic Dictionary Research Institute	日本電子化辞書研究所	
Japanese-to-English transfer dictionary	日英対照辞書	
basic lexicon	基本語彙	
co-occurrence dictionary	共起辞書	
collocation dictionary	共起辞書	
concept dictionary	概念辞書	
dictionary	辞書	noun
dictionary form	終止形	noun
generative lexicon	生成的辞書	
idiom dictionary	慣用語句辞書	
system dictionary	システム辞書	
terminology dictionary	専門辞書	
user dictionary	利用者辞書	
word collocation dictionary	単語共起辞書	

**Table 3. Sample Entries from lingdic.enja.utx**

#UTX-S 0.92; en-US/ja-JP; 2008-05-29; copyright: AAMT (2008); license: CC-by 3.0;		
#src	tgt	src:pos
head word	見出し語	noun
multiple-word string	複数語文字列	noun
single-word string	単一語文字列	noun
canonical form	基準形	noun
convention	慣例	noun
function word	機能語	noun
superlative	最上級	noun
unmarked sentence order	マークなしの文の順序	noun
base form	原形	noun
non-base form	非原形	noun
formulate	定型化する	verb
rule of thumb	原則	noun
compound	複合語	noun
marker	マーカー	noun
upper and lower case	大文字と小文字	noun
upper- and lower-case character	大文字と小文字の文字	noun
separator	分離記号	noun
string	文字列	noun
uninflected	活用されていない	adjective
orthographic	綴りの	adjective
lexicographical	辞書編集上の	adjective

**Table 4. Sample Entries from System D**