# Minna no Hon'yaku: A website for hosting, archiving, and promoting translations

**Masao Utiyama†, Takeshi Abekawa‡, Eiichiro Sumita† and Kyo Kageura§**

† MASTAR Project, National Institute of Information and Communications Technology, Japan

{mutiyama,eiichiro.sumita}@nict.go.jp

‡ Center for Informatics of Association, National Institute of Informatics, Japan

abekawa@nii.ac.jp

§ Graduate School of Education, University of Tokyo, Japan

kyo@p.u-tokyo.ac.jp

## 1   Introduction

In accordance with the rapid growth in the amount of information available in a variety of languages on the Web, translation needs and activities are sharply increasing. Many NGOs and NPOs that have been involved in translation as part of their activities since the pre-internet era have stepped up their translation activities; groups that were not previously involved in translation have started translating information related to their activities; online multilingual civic news sites have flourished; and a growing number of individuals have started translating online documents of their choice, in part or in full. Software localisation has come to constitute an important element in the distribution of software, and has become publicly visible due especially to the growth of free software.

Several notable features can be identified in the recent development of translation activities. One is the growth in the translation of online texts. This trend has gained impetus with the recent increase in the number of documents available under Creative Commons licenses, which helps promote the translation of online texts under clear licenses. A second feature is that online translation is mostly done by volunteer translators, who do not necessarily have expertise in translation. A third is that these translators rely heavily on dictionaries and Google search when translating. A fourth is that most of these translators do not use translation-aid tools. Finally, a fifth feature is that translation by volunteers has become an important source of providing unbiased or alternative information in an era of media concentration.

Against this backdrop, we have developed an integrated translation aid and hosting website, "Minna no Hon'yaku" ("Translation of/for/by everyone"; henceforth MNH), and made it publicly available since April 2009 at http://trans-aid.jp. The underlying concept of MNH is to make humans more intelligent by making machines more intelligent, and vice versa. MNH has two aspects: (i) a translation aid aspect that includes mechanisms for aiding online translators and facilitating the development of online translator communities by promoting inter-translator communication and collaboration; and (ii) a public media aspect that includes publishing translated information on the site and archiving translations with their originals. Section 2 of this

Figure 1. The top page of MNH

paper introduces the translation hosting aspect of MNH, and section 3 details the built-in translation aid editor QRedit. In section 4, we explain the community building functions as well as the educational functions provided by MNH. Section 5 describes the system's reference sources, and section 6 introduces its current state of usage. Section 7 concludes the paper.

## 2   MNH as a translation hosting site

### 2.1   Viewing documents

MNH provides information to those who access the MNH site looking for information. This is the same as any other news service or blog site. The only difference is that the information provided in MNH consists of translated documents. Figure 1 shows the MNH top page, which appears when internet users access the site. Anybody can read the translated documents published on the MNH site.

The main column in the centre provides contains newly uploaded translations; the right-hand column contains translations of "today's featured article" and "today's featured picture" from the English Wikipedia site [1]. The left-hand column provides a variety of information for users, i.e. how to use the MNH site, experimental use of QRedit, a link to the registration page, a

Figure 2. A personal user space in MNH

login box for registered users, and a list of translators who have recently published translations on MNH.

To navigate users to relevant documents, four sets of tags are displayed above the three columns: tags assigned to newly published translations; thematic tags; area tags; and genre tags. These tags are assigned by translators when they publish their translations on MNH. MNH provides a fixed number of pre-defined thematic, area and genre tags. In addition to these tags, translators can create their own tags. Users can also assign bookmark tags to published translations, which facilitates user-oriented organisation of translations on MNH.

## 2.2 Hosting translations

As a translation hosting site, MNH enables users to edit, store and publish their translations on MNH through their web browsers. In order to use these functions, one needs to obtain an MNH user account, which is free and open to anybody. Users can choose to remain anonymous, as MNH allows users to use pseudonyms. Registered users are provided with their own personal space,

公開設定　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　×

✎ 公開情報を入力して下さい。　　みんなの翻訳における著作権の考え方

あなたが翻訳の対象とした文書(原文)は、原文著者が明示的に許可している場合を除いて私的な利用その他など、著作権法で認められている範囲でしか利用できません。
原文著者は、その翻訳を公開しても良いと(あなたやみんなに)許可をしていますか？

◉ はい ○ いいえ

| 使用許諾条件 | あなたの文書(訳文)の使用許諾条件を選択して下さい。 | |
|---|---|---|
| | 🗁 クリエイティブコモンズの場合 | |
| | ◉ 表示 | この作品を利用するときには原著作者のクレジットを表示して下さい(説明🔗) |
| | ○ 表示-継承 | 「表示」に加えて、2次的著作物をこの作品と同一の使用許諾条件として下さい(説明🔗) |
| | ○ 表示-非営利 | 「表示」に加えて、非営利での使用のみを許可します(説明🔗) |
| | ○ 表示-非営利-継承 | 「表示-継承」に加えて、非営利での使用のみを許可します(説明🔗) |
| | ○ 制限なしに利用可能(パブリック・ドメイン) | パブリック・ドメインに属するものとして、世の中に寄贈します(説明🔗) |
| | 🗁 その他の場合 | |
| | ○ GNU Free Documentation License | GNU Free Documentation License(説明🔗) |
| | ○ 上記以外 | 「あなたの文書から2次的著作物を作成し、それを公開しても良い」という条件に矛盾しない上記以外の使用許諾条件を記述してください |
| アクセス権限 | もし他の人が、あなたの文書を翻訳等に利用することを望まない場合には、一般への公開はしないで、非公開文書として下さい | |
| | ◉ 非公開 | あなただけがこのファイルを閲覧・編集できます |
| | ○ 公開 | みんながこのファイルを閲覧できますが、編集ができるのはあなただけです |
| | ○ 公開・編集 | あなたに加えて、指定された人やグループがこのファイルを自由に閲覧・編集できます |

閉じる

Figure 3. The popup window that asks users for the copyright status of a document

in the same manner as ordinary blog services. Users can edit their personal information; make directories to manage translations; and make use of a variety of translation-aid and community building functions, etc. Figure 2 shows an image of a user's personal space when logged in.

The main differences between MNH and blog sites are:

1. As MNH is designed to host hosting translated documents, the documents that registered users publish are translations. Users make translations using the integrated translation-aid environment QRedit, the details of which will be explained in the next section. When users save their translations with permission to publish, they are published on MNH.

2. Translations made by users are published on the MNH top page. Unlike blog sites, the personal pages of registered users are not explicitly designed for the publication of their translations, although anyone can access the personal page of a registered user by inputting their user name as follows: `http://trans-aid.jp/users/?user_id=USER_NAME`.
Further, unlike most blog sites, individual users cannot change the design of their own space.

The clarification of copyright is a critical issue in hosting and publishing translated documents. Though it is users who take responsibility for copyright issues, MNH introduces three steps to make sure that users respect copyright. First, in the process of registration, users must agree to terms of use that request them to abide by copyright laws. Second, when users open QRedit to translate a document, MNH asks them to confirm that they follow copyright laws. Unless users agree, they cannot open QRedit (though they can elect to make a general agreement the first time they use QRedit so that they can open it without this step from then on). Third, when users save a translation, the system requires the clarification of the copyright status of the original document. MNH first asks whether the author of the original document permits the publication of translations. If users answer "yes", then a popup window appears. Figure 3 shows the popup window. In this window, users specify the copyright status, among the alternatives given

by MNH: Creative Commons BY, BY=SA, BY=NC, BY=NC=SA, GNU Free Documentation Licence, and others. Users then need to choose whether they want to publish the translation or keep it private (there is also a third choice, which we will elaborate on in section 4).

Translations of documents whose copyright permits the publication of translations are made publicly available on the MNH top page (as shown in Figure 1) when translators choose to publish the translations, while other translations made using QRedit on MNH are restricted to personal use or to use among a limited group of users. We will come back to the mode of limited use in section 4, in relation to the community-building functions of MNH.

# 3    QRedit: An integrated translation-aid editor

## 3.1    Needs of online volunteer translators

Volunteer translators involved in translating online documents have varied backgrounds. Some are professional translators volunteering part-time, some are interested in the topic, and some translate as a part of their NGO activities. They nonetheless share several basic characteristics:

1. They are native speakers of the target language (TL), which is mostly the case for any sort of translator;

2. Many do not have a native-level command of the source language (SL);

3. Most do not use translation aid systems or MT, because (i) most of these are commercial products that are too expensive for volunteer translators, and (ii) they are not readily accessible online;

4. They want to reduce the subjective burden involved in the process of translation. The reduction of time is also very important but perhaps less important for volunteer translators than for professional translators;

5. They spend a large amount of their translation time (about 30 to 60 per cent) looking up reference sources such as dictionaries or Google (cf. [2]);

6. The smallest basic unit of translation is perceived as a paragraph, not a sentence, and "at a glance" readability of the SL text is very important.

These traits define the basic features that any usable translation editor should provide for online volunteer translators:

1. It should be accessible online, through Web browsers;

2. It should provide high-quality dictionaries and an interface that promotes translators' awareness of dictionary entries for difficult expressions, so that they do not overlook, for example, idiomatic expressions [3];

3. It should provide a connection to Internet resources and should enable Internet searches to be initiated seamlessly from the translation editor environment;

4. It should provide a functional line of flow and interface that enables translators to maintain the rhythm of translation while using necessary functions such as dictionary lookup and Google search.
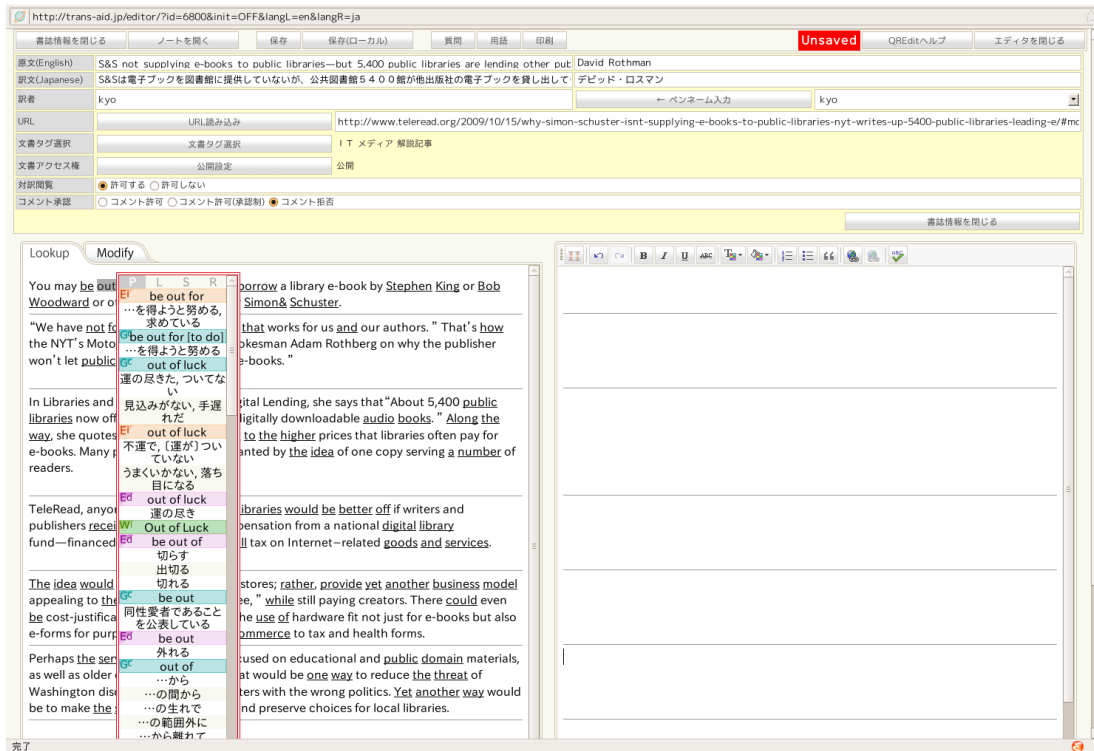
Figure 4. QRedit as it appears after start-up

## 3.2 Features of QRedit

QRedit has been developed as a translation-aid editor that provides various support functions to online translators based on the aforementioned requirements [3]. The user can open QRedit simply by clicking the QRedit start button, which appears immediately below the top banner of the MNH page when the user is logged in (see Figure 2). The user can either specify the URL of the document that s/he wishes to translate before clicking the start button or copy-and-paste the source document after opening QRedit. Figure 4 shows an image of QRedit after it has been opened. At this stage, (i) the bibliographic information area appears at the top of the window and (ii) dictionary lookup is activated in the source text area (which is on the left in Figure 4; users can rearrange the position of these areas). The bibliographic information area can be closed to enable the user to take full advantage of the screen for translating the document. Both the source and target text areas are divided into paragraph units. Scrolling of the source and target areas can be synchronised for both directions or from one to the other direction.

In the process of translation, translators can look up information by clicking words or phrases in the source text area. Figure 4 shows how dictionary lookup appears. When the user clicks a lexical token in the source text area, a small popup window appears, which displays translations from a high-quality dictionary provided by Sanseido [4], a few free dictionaries, and bilingual entries from Wikipedia (we will return to this point in section 5). Translation sentence pairs automatically collected from the Web and also provided by some NGO/NPO groups are also available for lookup [5], though not through the small popup window but by activating the sentence lookup session.

In addition to a simple dictionary lookup, QRedit also provides a flexible idiom/multi-word unit lookup mechanism. For instance, it can automatically look up the dictionary entry "with one's tongue in one's cheek" for the expression "He said that with his big fat tongue in his big fat cheek" or "head screwed on right" for "head screwed on wrong" [6].

QRedit provides a stratified reference lookup interface, which distinguishes three user awareness levels depending on the type and nature of the reference unit. These awareness levels are reflected in the way the reference units are displayed, for example, change in the background color or the use of an underline. The different awareness levels are assigned to reference units from a variety of reference sources, according to the criteria of "composition", "difficulty", "specialty" and "resource type" [7].

Starting from this small popup window, the user can trigger four types of action, corresponding to "P", "L", "S" and "R" shown on the top line of the popup window. In "P" mode (which is the default, though users can change their default settings), the user can paste any string in the popup window into the target text area by clicking it. If the user chooses "L" and then clicks the source language entry, a new, larger window appears, in which the full dictionary information is available. "S" mode enables the user to activate Google search by clicking a string in the small popup window. "R" activates the function of registering user defined translations.

Keyboard control remains fixed on the target text area, even when the user activates dictionary lookup by clicking source text tokens in the source text area using the mouse. This enables the user to maintain the rhythm of translation and focus on the making of target texts.

In an informal preliminary experiment, we observed two major effects of using QRedit:

1. Translation time is reduced by 20 to 30 per cent, mainly because QRedit dramatically reduces the time necessary for dictionary lookup and web search; and

2. The quality of draft translations is improved, because paragraph and textual reading is not hampered by a cumbersome process of reference lookup.

Though further detailed and systematic research is needed on this point, QRedit provides most of the translation-aid functions required by online volunteer translators.

# 4   Community building and implicit translation training

Although the main target of MNH is individual translators and groups of translators who already know one another (such as a team of volunteer translators working in an NGO), MNH has several Web 2.0 features to facilitate the community building among those who register at MNH. It has a social tagging system that allows users to tag translations with their favorite keywords; it allows users to submit questions about translation to MNH, so that users can help each other; and it also allows users to issue request for translation. Message exchange function is also provided in MNH, so that users can communicate with each other within MNH, even if they do not know each other's external e-mail addresses or contacts information.

Users can issue permission for other translators to edit their original translation. Such permission can be open-ended, or restricted to a particular group of users. This makes it easy for groups of translators sharing the same topics of interest to form loosely networked communities.

Figure 5. A comparative view of different versions of a translation

The most important contribution of this function to volunteer translators, however, is its implicit translation training function for inexperienced translators. Many NGOs and NPOs constantly face the problem of a lack of good volunteer translators as well as high turnover among those they do have. Due to the low retention rate of volunteer translators, the core, experienced translators become busier, leaving them no time to give advice to inexperienced translators, thus creating a vicious circle.

The implicit translation training in MNH works as follows:

1. An inexperienced translator, after translating a document, gives permission for experienced translators to edit the draft translation. Experienced translators edit the draft, and finalise the translation. Note that this is the normal task flow adopted in many NGOs/NPOs.

2. MNH keeps a log of these finalised translations (a maximum of 10 versions for any one document). Thus, a record of how experienced translators modified the draft translation is retained in the MNH database.

3. Translators who have the permission to edit can see the difference between any pair of translation versions kept in the MNH database. Inexperienced translators can therefore check where and how their translations were modified by experienced translators. Figure 5 shows the comparative view of two different versions of a translation.

Although this mechanism still requires that inexperienced translators make the effort to learn from the modification logs, it can potentially provide great assistance in improving their skills,

as they can take advantage of information that has not been accessible in the traditional working environment.

# 5 Reference resources

One of MNH's main features is that it provides rich – though limited in terms of language pairs – reference resources necessary for translation in QRedit. One of the most frequent requests made by volunteer translators is easy lookup of high-quality dictionaries and Wikipedia. Terminological management and TM are considered by many volunteer translators as desirable but not essential.

**Bilingual dictionaries:** The supply of a high-quality dictionary is essential to attract volunteer translators to MNH. We negotiated with one of the most prestigious dictionary publishers in Japan, Sanseido, and obtained permission to use *Sanseido's Grand Concise English-Japanese Dictionary*, which contains 360,000 entries. This is one of the most highly regarded English-Japanese dictionaries for translators. In addition, we provide the free Japanese-English dictionary EDICT [8] for both English to Japanese and Japanese to English translations. English-Chinese and Chinese dictionaries are to be supplied soon, which will extend MNH to Chinese.

**Encyclopedic information:** QRedit provides entries from Wikipedia. When there are corresponding Japanese and English entries, users can look up English to Japanese or Japanese to English translation information extracted from Wikipedia entries. When there are only English or Japanese entries, users can look up monolingual information.

**Terminology and proper names:** When MNH was made public, QRedit did not provide rich terminological resources. Since then, large terminologies of medicine and law have been uploaded by developers and users and made public. We have also developed a Web crawler to collect bilingual terminologies from the Web. This resource, consisting of about 2 million terms, will be made public as soon as information cleaning up is complete and copyright status is clarified [9]. An automatically collected compilation of proper-name translation pairs, containing around 400,000 proper-name translations, is also to be made available [10]. In addition, individual users can register terms either manually or using the term extraction and management mechanism, in which users can extract terms from their own or specified text pairs, check the result of extraction, and register them as their own terminological data.

**Translation memory:** The translations published on MNH are analysed and recycled as a TM resource. Major NGO users, such as Amnesty International and Democracy Now! Japan, have also provided their translations as TM for MNH. We also provide copyright free translation pairs collected from the Web (approximately 1.4 million English words) [11]. These are not currently retrievable from the QRedit small lookup window; another window must be used to retrieve TM. We will integrate TM lookup into QRedit as soon as possible.

# 6  Current state of usage

We made MNH publicly available on 7th April, 2009. Figure 6 shows statistics on MNH as of 20th October, 2009. The upper panel shows the number of registered users; the middle panel the number of registered documents (the top line shows the total number of documents, including both published and unpublished documents, and the lower line the number of published documents); and the bottom panel the number of registered bilingual term pairs (the top line shows the total number of terms, including both published and unpublished terms, and the lower line the number of published terms).

In the six months since MNH has become publicly available, it has attracted more than 900 registered users, of which 39 publish their translations through MNH. Amnesty International Japan and Democracy Now! Japan, as well as the Japanese translation team of GlobalVoices Online and some members of Translators United for Peace (TUP) are using MNH. A few universities use MNH in their lecture or seminar courses. One publisher, Sakuhinsha, has used MNH on an experimental basis to translate books into Japanese. One title has already been published [12], and others are due to follow.

Though the growth of registration has slowed down recently, we are expecting more users to register as we have just made the site fully compatible with Safari on the Mac and IE8 on Windows (previously, the full range of functions was only available with Firefox or Google Chrome) and are planning to launch a second-stage campaign for attracting NGOs/NPOs as well as publishers involved in translation work.

As for the number of translations, MNH currently has 2,340 translated documents, among which 1,235 have been published on the site. The other documents remain unpublished, or the translators may use them on their own sites, outside MNH. The total of 2,340 documents amounts to 17,066 translation pairs as counted by sentences. Since the end of May, the growth of translated documents has been linear.

Nearly 50,000 bilingual term pairs have been registered to MNH so far. Looking at the bottom panel of Figure 6, we can observe that most of these pairs were registered in April and May, 2009. This is because large terminological lists were uploaded during that period. Although the growth rate of bilingual terms registered to MNH has slowed down, it is still increasing steadily.

# 7  Conclusions

In this paper, we have introduced the basic features of the translation hosting site Minna no Hon'yaku (MNH) and QRedit, a translation-aid editor provided as part of MNH. The goal of MNH is:

1. to provide a large-scale archive of translations and their originals available to all under clear licenses, which constitutes a translation news site, and

2. to facilitate translation and nurture online translation communities

by providing online volunteer translators with a hosting service for translation activities and an easy-to-use translation-aid editor.
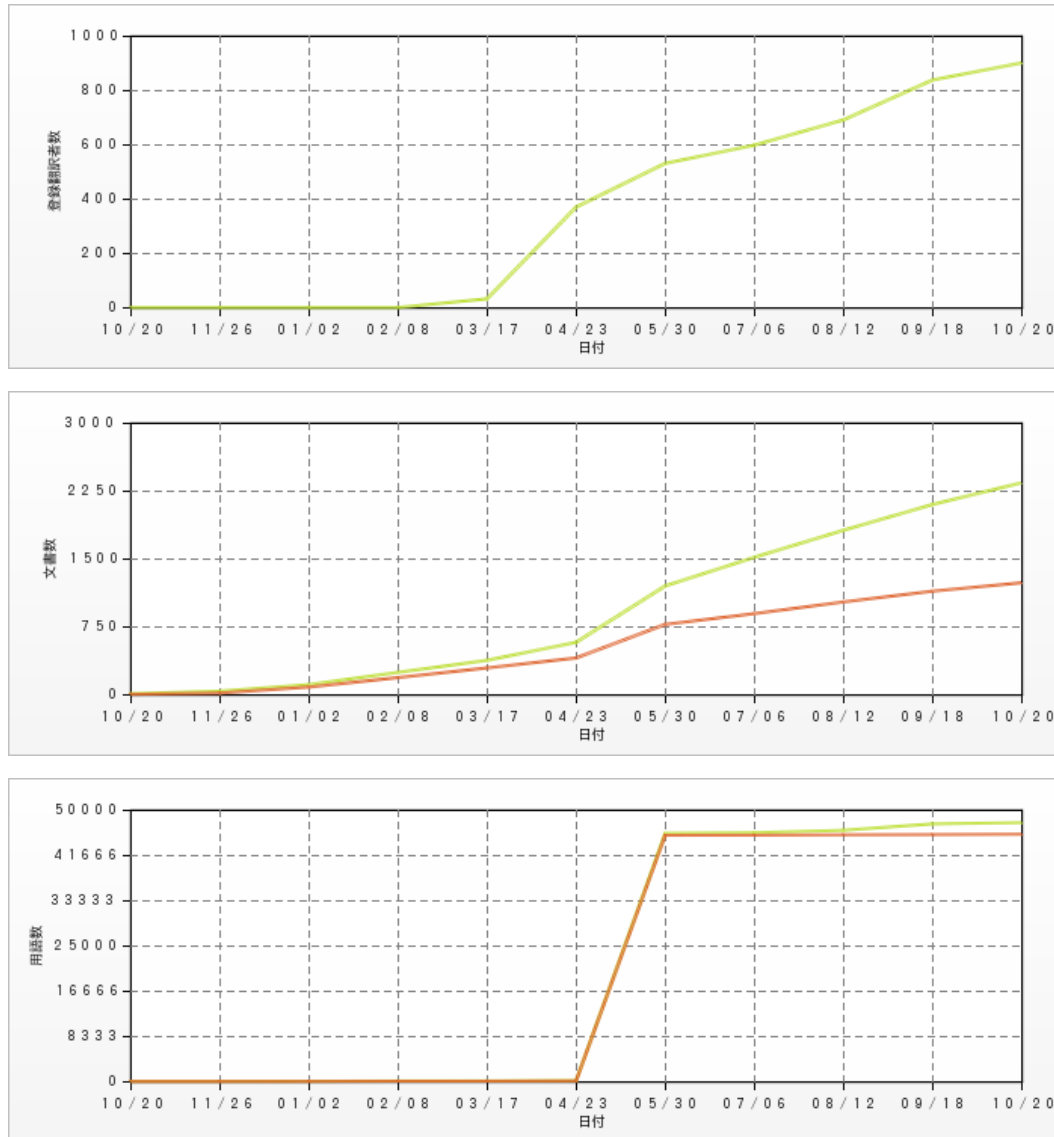
Figure 6. Some statistics on MNH

As mentioned above, according to an informal interview with three MNH users, translation time was reduced on average by 25 per cent, and the quality of the first, draft translation was improved by using MNH. Although we need to carry out systematic evaluation in order to confirm these reports, the fact that the number of users as well as the number of translated documents has steadily increased shows the practical usefulness of MNH. Though seemingly similar services have become available, such as Google Translator Toolkit [13], the translation-aid functions and hosting functions provided by MNH are unique and of higher quality. The basic concept of MNH is, as mentioned in the introduction, to make humans more intelligent by making computers more intelligent, and vice versa. Our emphasis at present is on human translators, but the translations accumulated in MNH are to be used to improve MT and the translation-aid functions of MNH further.

The biggest limitation of MNH for now is that it can only deal with English to Japanese and Japanese to English translations. This limitation of language pairs is inevitable in the short term, because to materialise fully the design concept of MNH, we need high-quality dictionary resources, such as the Sanseido dictionary we are currently using [4]. In the second stage, we will extend MNH to English to Chinese and Chinese to English translation. We hope that by extending the system to cover other languages with the cooperation of dictionary publishers, MNH will provide a social infrastructure for truly universal multilingual communication in the not-so-distant future.

## Acknowledgements

## References

[1] Wikipedia. `http://en.wikipedia.org/`

[2] Alain Désilet, Louise Brunette, Christiane Melanon and Geneviève Patenaude (2008). "Reliable innovation: A tecchie's travels in the land of translators," *8th AMTA Conference*, pp. 339–345.

[3] Takeshi Abekawa and Kyo Kageura (2007). "QRedit: An integrated editor system to support online volunteer translators," *Digital Humanities 2007*, pp. 3–5.

[4] Sanseido (2006). *Grand Concise English-Japanese Dictionary*. Tokyo: Sanseido.

[5] Masao Utiyama and Mayumi Takahashi (2003). *English-Japanese translation alignment data.* `http://www2.nict.go.jp/x/x161/members/mutiyama/align/index.html`

[6] Koichi Takeuchi, Takashi Kanehila, Kazuki Hilao, Takeshi Abekawa and Kyo Kageura (2007). "Flexible automatic look-up of English idiom entries in dictionaries," *Machine Translation Summit XI Proceedings*, pp. 451–458.

[7] Takeshi Abekawa and Kyo Kageura (2007). "A translation aid system with a stratified lookup interface," *Proceedings of the 45th ACL Annual Meeting Demos and Poster Sessions*, pp. 5–8.

[8] Jim Breen. EDICT. `http://www.csse.monash.edu.au/~jwb/edict.html`

[9] Takeshi Abekawa and Kyo Kageura (2009). "QRpotato: An exhaustive collector of bilingual terms from the web," *Proceedings of the 15th Annual Meeting of the Japan Society for Natural Language Processing* (in Japanese).

[10] Satoshi Sato (2009). "Crawling English-Japanese person name transliterations from the Web," *WWW 2009*, p. 1151–1152.

[11] Masao Utiyama and Mayumi Takahashi (2003). "English-Japanese translation alignment data," `http://www2.nict.go.jp/x/x161/members/mutiyama/align/index.html`

[12] Sakuhinsha. `http://tssplaza.co.jp/sakuhinsha/`

[13] `http://translate.google.com/toolkit/`