

Un paso adelante

Plan de tecnologías lingüísticas libres

Juan Rafael Fernández García

Copyright © 2003 Juan Rafael Fernández García

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.1 or any later version published by the Free Software Foundation.

\$Id: un-paso-adelante.xml,v 1.6 2004/07/29 09:11:43 olea Exp \$

Historial de revisiones

Revisión 0.9	2003-09-30	jrf
‘Demostración’ de las insuficiencias de .po		
Revisión 0.8	2003-09-29	jrf
Versión ‘La luz del sol ilumina las ideas’		
Revisión 0.7	2003-09-07	jrf
Versión enviada al VI Congreso Hispalinux		
Revisión 0.6	2003-03-17	jrf
Superada tortura de xmllint, victorioso		
Revisión 0.5	2003-02-26	jrf
La información pasa a doc-traduccion-libre		
Revisión 0.4	2003-02-25	jrf
Modificaciones menores del marcado		
Revisión 0.3	2003-02-24	jrf
Revisión del xml		
Revisión 0.2	2003-02-20	io
Primera versión xml		
Revisión 0.1	2002-17-12	jrf
Versión inicial (txt)		

Notas previas para una especificación de requisitos de las herramientas lingüísticas de TLDP (La comprensión plena de este documento requiere familiaridad con las herramientas lingüísticas o la lectura previa de doc-traduccion-libre — <http://es.tldp.org/Articulos/0000otras/doc-traduccion-libre/> [<http://es.tldp.org/Articulos/0000otras/doc-traduccion-libre/>], documento producido en el seno de *TLDP-ES*)

En este artículo el autor presenta y justifica sus propuestas con respecto a la adopción de estándares libres de marcado e intercambio de información lingüística y el desarrollo de herramientas libres de ayuda a la escritura y traducción, de manera que puedan someterse a discusión pública y adoptarse como especificaciones de TLDP. Es en cierto modo un apéndice ejecutivo, las conclusiones de la exposición realizada en el documento anterior.

- Última versión: <http://es.tldp.org/especificaciones/herramientas-linguisticas/herramientas-linguisticas/>
[<http://es.tldp.org/especificaciones/herramientas-linguisticas/herramientas-linguisticas/>]
- Fuente: <http://cvs.hispalinux.es/cgi-bin/cvsweb/esp-herramientas-linguisticas/un-paso-adelante.xml>
[<http://cvs.hispalinux.es/cgi-bin/cvsweb/esp-herramientas-linguisticas/un-paso-adelante.xml>]

Tabla de contenidos

Introducción	
¿Por qué?	
Propuesta	
Un proyecto paralelo: Okapi Framework	
Al grano	
Propuesta: Líneas generales	
Diccionarios desde el punto de vista lexicográfico	
Correctores ortográficos	
Correctores ortotipográficos	
Terminologías y herramientas de gestión terminológica (TMS)	
Correctores gramaticales y de estilo	
Corpora monolingües	
Corpora paralelos y multilingües	
Memorias de traducción	

Introducción

La temática de este documento es lo que en la jerga de la industria se llama «*localización*» o incluso «*globalización*», directamente relacionadas con la traducción y sus herramientas. Me he esforzado en atenerme a desarrollos posibles a corto y medio plazo, pero continuamente me ha asaltado la tentación de entrar en el tema de la traducción automática. Queremos, a corto plazo, que las máquinas nos ayuden a traducir; a medio plazo el objetivo es que traduzcan por nosotros. Lo que pasa es que los dos objetivos convergen en dotar de inteligencia a las máquinas. Pienso que cualquier asistente a la traducción necesita reglas: morfológicas, sintácticas, semánticas, necesita una gramática y un «conocimiento del mundo». ¿Cómo traduciría un programa sin inteligencia oraciones como ‘spirits sold here’? ¿cómo puede un humilde corrector ortográfico distinguir ‘e implementado’ de ‘he implementado’?

Es la hora de que el software libre llegue a la madurez. Eso implica también que conozcamos las tendencias de la industria y adoptemos y trabajemos con los estándares. Este documento trata de señalar el camino para el desarrollo de nuestras herramientas de CAT, con el ojo puesto en el desarrollo futuro de herramientas de traducción más o menos automática. Ahora hay que ponerse técnicos y me temo que la lectura no va a ser fácil; no soy un experto y escribo mientras aprendo así que es probable que se encuentren inexactitudes que será conveniente corregir. De todos modos no pretendo más de plantear unas perspectivas y unas posibles líneas de desarrollo; está claro que LuCAS y TLDP son ejemplos de éxito como modelos de cooperación libre y que el actual proceso de creación de la editorial libre y la conversión de la documentación a DocBook XML son grandes ideas.

Dice Ismael Olea, hablando del ciclo de desarrollo de TLDP,

Hay dos niveles de objetivos a cumplir:

- herramientas para el fácil mantenimiento y traducción de nuestra documentación
- herramientas lingüísticas de calidad

Mi objetivo es avanzar en la segunda opción.

¿Por qué?

Pensando sobre el esquema [<http://olea.org/conferencias/doc-conf-creando-lucasv4/herramientas-reproduccion-publicacion.png>] de Ismael Olea ‘Arquitectura para reproducción y publicación’ del borrador de “ LuCAS v4: creando el sistema de publicación de TLDP-ES ”, creo que es el momento —intentaré explicar por qué— de integrar en el esquema las herramientas lingüísticas y dar el paso de adoptar los estándares y potencialidades de la ingeniería lingüística actual.

Dependiendo de cómo nos planteemos la realidad de nuestro fondo documental, podemos pensar en él como en una gran biblioteca electrónica, como un repositorio digital (como el DSpace del MIT) o como un gran corpus lingüístico. Cada punto de vista plantea posibilidades distintas: en tanto que biblioteca, debemos escuchar las aportaciones de los bibliotecólogos. Debemos también acordar los metadata y el formato de registro de los documentos, para que las búsquedas sean cada vez más eficaces. Yo, como lingüista, propongo que lo miremos con un nuevo enfoque: somos un gran equipo de traductores y redactores de documentación. Una de nuestras prioridades debe ser mejorar nuestras herramientas de ayuda a la escritura.

Disponemos de un fondo de varios cientos de megas de documentación técnica, en proceso de transcripción a DocBook XML, en su mayor parte traducciones de obras o documentación de tema informático que también son de licencia libre y disponemos además de las traducciones a multitud de idiomas. Ninguna empresa —¿quizás alguna universidad?— dispone de un fondo documental como el nuestro para la elaboración de terminologías específicas, para la confección de un corpus tan completo de la materia ni para la creación de memorias de traducción. Ninguna empresa dispone del número de traductores de que dispone la comunidad GNU Linux. Y sin embargo las características de esta comunidad (en su mayor parte estudiantes y profesionales de la informática) hace que el modo de trabajar y las herramientas estén a años luz de las posibilidades que presenta la investigación actual. No se conocen las experiencias que se realizan en la Unión Europea ni las herramientas que se utilizan, ignoramos los proyectos universitarios y los intentos de estandarización que está realizando la industria.

Y la distancia con respecto a la evolución del software propietario, lamentablemente, se amplía en lugar de acortarse. La comunidad, bajo el nombre de HispaLinux, o de TLDP o de FSF o del que sea, debe funcionar como una entidad y personarse y participar en los foros donde se están creando estándares (LISA, ELRA...) y en los proyectos financiados por la Unión Europea (debemos aspirar a recibir financiación de las instituciones y de las empresas). Es la hora de preguntarse si el modelo de voluntarios a tiempo parcial es eficaz para los objetivos que nos proponemos.

La crisis: más allá de los ficheros .po

No podemos olvidar que gettext y .po resuelven algunos de los problemas con los que tiene que enfrentarse todo software de traducción: la extracción de las cadenas que deben ser traducidas, la segmentación y la referencia al lugar en el código al que pertenece la cadena, la alineación entre texto fuente y texto traducido, la capacidad de reutilizar traducciones cuando el fichero original ha cambiado, marcando las cadenas equivalentes y aquellas en las que la traducción es sólo aproximada (*fuzzy*) o no existe, y el informe automático del número de cadenas traducidas...

Existen varios intentos de ampliar el uso de nuestras herramientas más allá de la traducción de interfaces de usuario, a textos libres (ficheros no .po, XML especialmente), mediante el procedimiento de convertir de alguna manera los documentos a ficheros .po:

- *doc-i18n-tool* ((<http://mail.gnome.org/archives/gnome-doc-list/2001-October/msg00034.html>
[<http://mail.gnome.org/archives/gnome-doc-list/2001-October/msg00034.html>])
- *poxml* (de KDE)

- *po-debiandoc*

- *po-pod*, nuevo miembro de la familia de utilidades po-for-everything (po4a), tal como se presenta en la Debian Weekly News del 3 de dic. de 2002:

“ The goal of po-pod is to allow translators to work only with well known po files when translating pod documentation. The goal of po4a is to ease translations (and more interestingly, the maintainance of translations) by using gettext tools on areas where they were not yet expected. ” (<http://lists.debian.org/debian-i18n-0211/msg00009.html> [<http://lists.debian.org/debian-i18n-0211/msg00009.html>])

La referencia es <http://www.ens-lyon.fr/~mquinson/deb.html#po-pod> [<http://www.ens-lyon.fr/~mquinson/deb.html#po-pod>].

Segunda visita a KDE

La segunda y más novedosa parte del «*The KDE Translation HOWTO*» es el capítulo titulado «*Doc Translation*».

¹

To make the documentation easier to maintain, the source files ² were converted to PO format. This format was already used by GUI translators with great success. Matthias Kiefer and his contributors extended KBabel to accomodate for this new task (e.g. colored diff mode, translation database, extended line feed handling). So KBabel is probably a must have now not only for GUI translators but also for people working on doc translation.

Stephan Kulow enabled the KDE Help Center to parse XML(TM) files directly and to generate the HTML on the fly.

Los traductores del proyecto tienen la ayuda de las herramientas proporcionadas por el paquete *poxml*: *fixsgml*, *xml2pot*, *split2po* y por *xmlizer* y *checkXML*, de *kdelibs3-bin*.

Propuesta

Llevo tiempo preguntándome si estos tímidos intentos van bien encaminados. Mis objeciones son las siguientes:

1. *.po no es un formato bien marcado*

Me refiero, cualquiera que haya luchado con el Robot del Proyecto de Traducción Libre me entenderá, a que quede claro que en la tercera línea debe ir el *e-mail* del traductor y el *copyright*, por ejemplo. La DTD de *.po* está implícita de forma tácita en las herramientas disponibles, pero no se muestra en el formato.

2. *Me temo que .po no permite granularidad en la descomposición del párrafo/cadena en sus componentes.*

El formato nació para traducir *mensajes* (el trabajo del traductor consiste en generar las *msgstr* correspondientes a los *msgid* extraídos del programa), es decir, cadenas en una interfaz de usuario, sean los componentes de un menú o los anuncios que se le hace al usuario (‘el fichero ha sido guardado’). Pero es que ni siquiera para esta función es eficaz: la *segmentación* de la salida de *-h* de cualquier programa medianamente complejo (pienso ahora en las varias pantallas de *nmap*) es nula. ¿Cómo va a reutilizarse esta traducción?

La segmentación a nivel de párrafo en un texto XML es totalmente insuficiente: un párrafo puede ocupar varias páginas.

¹ <http://i18n.kde.org/translation-howto/doc-translation.html> [<http://i18n.kde.org/translation-howto/doc-translation.html>].

² El formato de la documentación había empezado siendo *.html*, fue LinuxDoc *.sgml* con KDE 1.x y DocBook *.sgml* en KDE 2.0.

3. En consecuencia de las insuficiencias anteriores, *no es posible una paralelización de calidad entre la unidad de traducción fuente y la unidad destino*, puesto que la única segmentación es en párrafos.
4. Sospecho que *cuando reconoce una cadena como «difusa» (fuzzy) está actuando únicamente al nivel de diff*.
No se utilizan las posibilidades de las expresiones regulares y de los algoritmos avanzados en el reconocimiento de patrones; no se avanza en la reutilizabilidad de las traducciones.
5. Aunque gettext es extremadamente eficiente en la extracción de cadenas cuando se aplica a un programa adaptado, *no hay modo de separar en un texto en XML las cadenas que deben ser traducidas de las que no es necesario o no deben traducirse*. No hay modo de utilizar los recursos de XLIFF.
6. Una última consecuencia de la no granularidad de la segmentación y alineamiento de los textos XML es la dificultad de crear memorias de traducción útiles.

Aunque a efectos del resultado sería indiferente el formato que utilizara internamente una utilidad de traducción, en nuestro análisis debemos considerar si el formato elegido es adecuado para nuestras necesidades (ya no la traducción de cadenas en tanto que mensajes de un programa informático, sino de documentación).

Quizás el camino no esté en convertir los ficheros xml en ficheros .po para poder traducirlos, sino por el contrario mantener gettext y el proceso de extracción de mensajes, que funciona bien y es estable, convertirlos en ficheros XML bien formados, y aplicar herramientas generalizadas que puedan utilizarse sobre cualquier fichero XML. Lo que propondría sería una especie de .po con marcas, utilizable además para texto libre.

Un proyecto paralelo: Okapi Framework

El proyecto del Marco de trabajo Okapi (<http://okapi.sourceforge.net/> [<http://okapi.sourceforge.net/>]) tiene licencia modelo MIT.

The goal of the «*Okapi Framework*» is to provide public specifications, open components and libraries to allow tools developers and localizers to build their own processes the way they want it, while staying compatible and interoperable with each others.

You can see the Okapi Framework as a set of links, of construction blocks, the glue you can use to put together a process that takes advantage of existing tools as you see fit, by allowing them to work more efficiently together.

The framework uses and promotes open standards when they exist. For the aspects where open standards are not defined yet, the framework offers its own proposals in an effort to prompt and help the definition of the missing open specifications. The ultimate goal is always to adopt the industry standards when they are defined.

Propone XLIFF 1.0 para la extracción de texto, TMX 1.3 para intercambiar TMs, TBX para el intercambio de terminologías y OLIF para el de lexicones (glosarios para sistemas de MT).

Al grano

En las secciones que restan de este documento vamos a seguir el procedimiento de discutir brevemente las necesidades y herramientas libres conocidas y acabaremos por hacer una propuesta, bien de adopción de un estándar o desarrollo de una aplicación, bien de estudio previo a tomar una decisión.

Son de varios géneros las utilidades de ayuda a la escritura y al traductor: diccionarios, concordancias, memorias de traducción, analizadores de todas clases. No hay espacio en este artículo para volver a hacer lo que como decíamos en el resumen inicial ya está hecho: subyace a todas las referencias la familiaridad tácita con las herramientas o con <http://es.tldp.org/Articulos/0000otras/doc-traduccion-libre/>

[<http://es.tldp.org/Articulos/0000otras/doc-traduccion-libre/>]. Es el momento de una revisión detallada y de tomar decisiones.

Propuesta: Líneas generales

En TLDP se usarán estándares libres y abiertos: DocBook, TEI, XCES... y cuando se utilice uno propio se crearán mecanismos para la importación/conversión (recordemos por ejemplo la utilidad incorporada a Mimersbrunn como conversor de .po a TMX).

El formato de los documentos será preferentemente DocBook XML, porque permite la separación lógica entre los niveles semántico y de presentación y da pie a desarrollos relacionados con la web semántica, y las herramientas trabajarán sobre XML.

La codificación UTF-8 es lógica en un proyecto con ambiciones de universalidad.

Diccionarios desde el punto de vista lexicográfico

En primer lugar evitemos una confusión generalizada, *de verba non est disputandum*. Diccionarios son listas de palabras de un idioma, entradas léxicas con su definición, tablas de equivalencias entre dos o más lenguas y cualesquiera otras variaciones que se nos planteen; nomenclatura, glosario, lexicón, lemario etc. son términos que se usan de manera no estable en la literatura y no deben impedir que nos entendamos³. Mucho más interesante es plantearnos el campo de realidad cubierto por el diccionario y la información que proporciona (morfológica, de uso, terminológica...). Por último, distinguiremos diccionarios destinados a ser usados por humanos (en forma de libro o como consultas a través de una interfaz) de aquellos creados para ser usados por máquinas. También deberemos distinguir los que tienen una ambición descriptiva y verbal de los diccionarios terminológicos, que son como sabemos normativos y su objeto son los conceptos de un campo. Al fin y al cabo por lo pronto sólo queremos traducir documentación informática ¿o no?

Discusión

Estamos en LuCAS-desarrollo reimpulsando nuestro propio diccionario inglés-español de informática Gaiat⁴ Por otro lado estamos intentando crear el mayor lemario del español⁵. Y periódicamente nos planteamos qué sentido tendría reclamar la liberación del diccionario de la RAE (probablemente muy escaso). Pero todos estos son pequeños pasos. Ponernos a escribir diccionarios nos convierte en lexicógrafos; pero escribir un diccionario técnico nos obliga a plantearnos los problemas de la terminología. Además, si logramos crear una estructura para la confección de terminologías (partiendo de nuestro corpus de informática pero no quedándonos ahí) podremos crear glosarios y diccionarios, presentar colocaciones y listas de ejemplos. Y hacer un uso normativo (estandarizado, coherente) de nuestra terminología. Estaremos por ejemplo en condiciones de definir como término informático el que aparece en nuestro fondo documental y no aparece recogido en nuestro diccionario general de español (por ahora inexistente).

Pero vayamos poco a poco: un diccionario pensado para su impresión en papel o en pantalla, destinado a ser consultado por humanos, no es lo mismo que una base de datos terminológicos (un lexicón computacional), pensada para su utilización en la traducción automática.

Sobre TEI P4

La versión XML de las *TEI P4 Guidelines*⁶, en su capítulo 13 'Terminological Databases', advierte

³Basta consultar la clasificación de dos páginas del término «diccionario» que hace Martínez De Sousa *Diccionario de lexicografía práctica* para llegar a la conclusión de que lo importante no es cómo se llama sino qué clase de diccionario queremos.

⁴ <http://cvs.hispalinux.es/cgi-bin/cvsweb/rl-dicc/> [<http://cvs.hispalinux.es/cgi-bin/cvsweb/rl-dicc/>].

⁵ <http://lemarios.olea.org> [<http://lemarios.olea.org>]

⁶ (<http://www.tei-c.org/Guidelines2/index.html> [<http://www.tei-c.org/Guidelines2/index.html>]). El documento puede descargarse comprimido como <http://www.tei-c.org/Guidelines2/p4html.tar.gz> [<http://www.tei-c.org/Guidelines2/p4html.tar.gz>]

Since its first publication, this chapter has been rendered obsolete in several respects, chiefly as a result of the publication of ISO 12200, and a variant of it (TBX) which has been recently adopted by LISA. Work is currently ongoing in the ISO community to define a generic platform for terminological markup (ISO CD 16642, TMF: Terminological Markup Framework), in the light of which it is anticipated that the recommendations of the present chapter will be substantially revised.

Sí nos interesa el capítulo 12 ‘Print Dictionaries’.

A simple dictionary entry may contain information about the form of the word treated, its grammatical characterization, its definition, synonyms, or translation equivalents, its etymology, cross-references to other entries, usage information, and examples.

Interesante parece estudiar como ejemplo de aplicación de TEI <http://www.human.toyogakuen-u.ac.jp/~acmuller/articles/ddb-ebti2001.htm> [<http://www.human.toyogakuen-u.ac.jp/~acmuller/articles/ddb-ebti2001.htm>]

Sobre WordNet

Información en Princeton WordNet (<http://www.cogsci.princeton.edu/~wn/w3wn.html> [<http://www.cogsci.princeton.edu/~wn/w3wn.html>]
ó

man wndb

WordNet es “an online lexical reference system whose design is inspired by current psycholinguistic theories of human lexical memory”.

EuroWordNet (<http://www.illc.uva.nl/EuroWordNet/> [<http://www.illc.uva.nl/EuroWordNet/>]). EuroWordNet was a European resources and development project (LE-2 4003 & LE-4 8328) supported by the Human Language Technology sector of the Telematics Applications Programme, 1996-1999

Problema

Me temo que la versión europea de wn tiene carácter no libre y sus herramientas son no libres (editor Polaris; el visor Periscope solamente es freeware)

Propuesta

1. El formato estándar para la distribución de diccionarios es .dict (un éxito, multiplicado casi por cuatro desde la primera redacción de doc-traduccion-libre [<http://es.tldp.org/Articulos/0000otras/doc-traduccion-libre/>])
2. Para consultas el estándar es la red dict, mediante clientes específicos o interfaz *web*
3. Otro tema es cómo escribir diccionarios: propongo con dudas adoptar TEI P4 como la DTD de nuestros diccionarios específicos. En este sentido, es necesario también completar el desarrollo de Giait y desarrollar herramientas generales que puedan utilizarse para el desarrollo de nuevos diccionarios
4. La gran pregunta, que normalmente no se hace cuando se está haciendo un diccionario, es qué información debe contener. Si queremos poder utilizar herramientas avanzadas de ayuda a la traducción nuestro diccionario debe tener información gramatical leíble por máquinas. Propongo el uso de XCES.

5. Una cuestión previa imprescindible es analizar ISO 12620 (sobre *recordable properties of terms*) y llegar a un acuerdo (parts of speech, gender, context, subject field...) sobre las propiedades que nos es útil recoger en el diccionario.

6. Nuestra urgencia real es crear terminologías unificadas.

7. Cómo incorporar términos a los diccionarios: ¿modelo ORCA? (contribuciones públicas, revisadas por un moderador) ¿modelo lista spanglish? (discusiones inter pares)

Utilizar herramientas de confección de glosarios. Debemos ir hacia diccionarios basados en *corpus*, con herramientas de extracción terminológica que cubren de forma exhaustiva un campo (y nos señalan fehacientemente qué términos faltan por definir) y abandonar el método manual de adición de entradas.

8. Cómo garantizar la calidad de las incorporaciones: Crear sistema de mantenimiento de calidad (incluir entre las marcas de cada término autor, fecha, revisión...)

Ismael Olea ha propuesto un sistema de ponderación de la autoridad de las aportaciones similar al del sistema *advogato* (<http://advogato.org/trust-metric.html> [<http://advogato.org/trust-metric.html>]).

Será necesario adaptar TEItools (<http://xtalk.msk.su/SGML/TEItools/index-en.html> [<http://xtalk.msk.su/SGML/TEItools/index-en.html>]).

Un caso práctico

Las dos primeras entradas de la letra «a» en las fuentes del diccionario Gaiit son

```
@@ING A, Ampere, Amps
@@CAS Ampère (amperio), amperios.
@@DIC Vocablo definido únicamente para el Diccionario
@@FIN
```

```
@@ING A Programming Language", APL
@@CAS Lenguaje de Programación APL
@@GLO APL es la sigla de "A Programming Language" (Un Lenguaje de
@@GLO Programación), que fue un libro escrito en 1962 por el creador del
@@GLO lenguaje, Kenneth E. Iverson. Basado en lo que antes se había conocido
@@GLO como la Nomenclatura Iverson, el APL es un lenguaje de programación
@@GLO extremadamente conciso, diseñado para el manejo de los arreglos
@@GLO (arrays). Los arreglos pueden ser escalares, vectoriales, tablas o
@@GLO matrices de dos o más dimensiones, pudiendo estar compuestos de
@@GLO información numérica o alfanumérica. Bajo la conducción de Iverson,
@@GLO IBM introdujo en 1966 el APL\360.
@@GLO Como el APL evita la introducción de las computadoras personales,
@@GLO originalmente se lo empleó únicamente en mainframes. Desde 1983, han
@@GLO estado disponibles las versiones de APL para las PC. Debido a su
@@GLO conjunto de caracteres especiales expandidos, el APL requiere un
@@GLO teclado especial, o el uso de macros, para el ingreso de datos.
@@GLO Las versiones actuales de APL para mainframes y PC se denominan
@@GLO APL2.
@@GLO
@@FIN
```

Convertidas a formato .dict (descomprimido) actualmente quedan así:

A, Ampere, Amps Ampère (amperio), amperios.

A Programming Language", APL Lenguaje de Programación APL

Vemos que, aparte de errores tipográficos que habrá que corregir, se ha perdido la información de glosario y el marcado (lo que es inglés, castellano, comentario, glosario...)

¿Cómo aparece?

```
11262928 23 n 03 ampere 0 amp 0 A 0 003 @ 11259168 n
0000 #p 11263273 n 0000 %p 11263164 n 0000 | the basic unit of
electric current adopted under the System International d'Unites;
"a typical household circuit carries 15 to 50 amps"
```

¿Cómo quedaría todo esto en TEI P4? Crearemos el fichero ejemplo.tei.xml (el siguiente fragmento de texto está en UTF-8)

```
<!-- %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% empieza ejemplo.tei.xml -->
<?xml version="1.0" encoding="UTF-8" ?>
<!DOCTYPE TEI.2 PUBLIC "-//TEI P4//DTD Main Document Type//EN"
  "/usr/local/share/sgml/tei/dtd/tei2.dtd" [
  <!ENTITY % TEI.XML          'INCLUDE' >
  <!ENTITY % TEI.dictionaries 'INCLUDE' >
]>

<body>
  <div0 type='dictionary'>
    <!-- English-español -->
    <entry>
      <form type="abbrev">
        <orth>A</orth>
      </form>
      <form type="full">
        <lbl>abbreviation for</lbl>
        <orth>Amper</orth>
        <trans>
          <tr>A</tr>
        </trans>
      </form>
    </entry>

    <entry>
      <form>
        <orth>Amper</orth>
      </form>
      <gramGrp>
        <pos>n</pos>
      </gramGrp>
      <def></def>
      <trans>
        <tr>Amperio</tr>
      </trans>
      <gen>m</gen>
    </entry>

    <entry>
```

```

    <form type="abbrev">
<orth>APL</orth>
    </form>
    <form type="full">
      <lbl>abbreviation for</lbl>
      <orth>A Programming Language</orth>
      <trans>
        <tr>APL</tr>
      </trans>
    </form>
  </entry>

  <entry>
    <form>
      <orth>A Programming Language</orth>
    </form>
    <gramGrp>
      <pos>n</pos>
    </gramGrp>
    <def>
      APL es la sigla de "A Programming Language" (Un Lenguaje de
      Programaci3n), que fue un libro escrito en 1962 por el creador del
      lenguaje, Kenneth E. Iverson. Basado en lo que antes se habAa conocido
      como la Nomenclatura Iverson, el APL es un lenguaje de programaci3n
      extremadamente conciso, diseAado para el manejo de los arreglos
      (arrays). Los arreglos pueden ser escalares, vectoriales, tablas o
      matrices de dos o mAs dimensiones, pudiendo estar compuestos de
      informaci3n numArica o alfanumArica. Bajo la conducci3n de Iverson,
      IBM introdujo en 1966 el APL\360.
      Como el APL evita la introducci3n de las computadoras personales,
      originalmente se lo emple3 unicamente en mainframes. Desde 1983, han
      estado disponibles las versiones de APL para las PC. Debido a su
      conjunto de caracteres especiales expandidos, el APL requiere un
      teclado especial, o el uso de macros, para el ingreso de datos.
      Las versiones actuales de APL para mainframes y PC se denominan
      APL2.
    </def>
    <trans>
      <tr>Lenguaje de Programaci3n APL</tr>
    <gen>m</gen>
    </trans>
  </entry>
</div0>
</body>
<!-- %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% acaba ejemplo.tei.xml -->

```

¿Cómo convertir esta fuente en un documento distribuible? Utilizaremos lib-saxon-java, passivetex y las hojas de estilo XSL creadas por Sebastian Rahtz (<http://www.tei-c.org/Stylesheets/teixsl.html> [<http://www.tei-c.org/Stylesheets/teixsl.html>]). Lo siguiente depende de la instalación de cada uno; yo utilizo en estos momentos una Debian Sarge. Passivetex es un paquete *apt-gettable*; descargo las DTD's TEI y las hojas de estilo y las instalo en /usr/local/share/sgml/tei/

La conversi3n a fichero .pdf me da errores; te3ricamente se harAa

```

java -classpath /usr/share/java/saxon.jar \
com.icl.saxon.StyleSheet -o ejemplo.tei.fo ejemplo.tei.xml \

```

```
/usr/local/share/sgml/tei/xsl/fo/tei.xsl
```

```
pdfxslt ejemplo.tei.fo
```

```
pdfxslt ejemplo.tei.fo
```

Para generar el fichero `ejemplo.resultado.html`

```
java -classpath /usr/share/java/saxon.jar \  
com.icl.saxon.StyleSheet ejemplo.tei.xml \  
/usr/local/share/sgml/tei/xsl/html/teihtml.xsl
```

Este es el fichero generado:

```
<!-- %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% empieza ejemplo.resultado.html -->  
<?xml version="1.0" encoding="utf-8"?>  
<!DOCTYPE html  
  PUBLIC "-//W3C//DTD HTML 4.0 Transitional//EN">  
  
<div class="teidiv">  
  <h2><a name="ej.utf8-div0-d0e3"></a>1.  
  </h2>
```

A

abbreviation for
Amp`re

A

Amp`re

n

Amperio
m

APL

abbreviation for
A Programming Language

APL

A Programming Language

n

APL es la sigla de "A Programming Language" (Un Lenguaje de Programaci3n), que fue un libro escrito en 1962 por el creador del lenguaje, Kenneth E. Iverson. Basado en lo que antes se haba conocido como la Nomenclatura Iverson, el APL es un lenguaje de programaci3n extremadamente conciso, diseado para el manejo de los arreglos (arrays). Los arreglos pueden ser escalares, vectoriales, tablas o matrices de dos o m3s dimensiones, pudiendo estar compuestos de informaci3n num3rica o alfanum3rica. Bajo la conducci3n de Iverson, IBM introdujo en 1966 el APL360. Como el APL evita la introducci3n de las computadoras personales, originalmente se lo emple3 unicamente en mainframes. Desde 1983, han estado disponibles las versiones de APL para las PC. Debido a su conjunto de caracteres especiales expandidos, el APL requiere un teclado especial, o el uso de macros, para el ingreso de datos. Las versiones actuales de APL para mainframes y PC se denominan APL2.

Lenguaje de Programaci3n APL
m

</div>
<!-- %%% acaba ejemplo.resultado.html -->

Correctores ortográficos

Tema relacionado con el de los diccionarios en tanto que listas de palabras.

«Nuestro» proyecto es el *COES*: Herramientas para Procesamiento de Lenguaje Natural en Espaol (<http://www.datsi.fi.upm.es/~coes/> [<http://www.datsi.fi.upm.es/~coes/>]), de Santiago Rodríguez & Jesús Carretero, que se distribuye como software de libre disposici3n desde finales de 1994.

COES consta de un diccionario de unos 53.000 t3rminos y un corrector ortogr3fico integrado en la utilidad Unix *ispell* y desarrollos derivados (*aspell*). Hay que sealar que se ha ampliado el conjunto de herramientas lingüísticas con un diccionario de sin3nimos/ant3nimos. Su particularidad es ser sensible a las reglas morfológicas de las palabras y no sólo a las raíces.

Un corrector ortográfico «inteligente» tiene que hacer algo más que comparar las palabras del texto con una lista de palabras correctas; para distinguir «a» de «ha» tiene que tener reglas — acabaremos necesitando un mínimo análisis morfológico y sintáctico.

Correctores ortotipográficos

Estudiar la implementación de OpenOffice.

Terminologías y herramientas de gestión terminológica (TMS)

Discusión

Las ventajas que ofrece la utilización de terminologías son evidentes: uniformidad entre las traducciones, corrección y reusabilidad. Debemos también en el seno de TLDP clarificar los problemas ante los que se enfrenta toda terminología:

1. eliminación de la sinonimia
2. reducción de la polisemia/homonimia
3. neologismos

Terminología y Memorias de traducción

Discutido en <http://www.lisa.org/sigs/phpBB/viewtopic.php?topic=22&forum=1&4> [<http://www.lisa.org/sigs/phpBB/viewtopic.php?topic=22&forum=1&4>]

Responde Kara Warburton:

There are many differences between Terminology and Translation Memory. The first that comes to mind is that terminology management must occur in the source language to ensure high quality source texts before a translation memory even exists. Inconsistency and inappropriate terms in the source language is an increasing and costly problem especially for global companies. And if the inconsistency is unchecked, then it breaks the TM — there will be no match with previous texts which are "supposed" to be the same. A second difference is that terminology management can document features of words in a much more granular fashion than translation memory which only records strings of text. Definitions, contexts, usage notes, references to related terms, and so forth can assist translators to select an appropriate target term when the TM does not suggest one. And if there is an exact TM match, what if the meaning is different which can frequently occur? Then the TM match is incorrect and the translator has to think more carefully about the terms. In this case, having a terminology system can provide the translator with the required information to correct the TM. In addition, grammatical information can be recorded which can provide data to more advanced systems down the road, such as search engines and machine translation engines.

Ingrid Haussteiner lo confirma:

I think your question is very interesting and I think from a translator's perspective this is not an either-or thing but rather a question of many ifs.

If you can choose, I would implement a terminology management system first, have people test and hone their discipline and raise their awareness of the difficulties involved in using consistent terminology (even more so if it is multilingual), analyzing choices and making suggestions (descriptive terminology work).

If you implement a translation memory system after some time, you might see that people slack off a bit as far as terminology work is concerned. The big argument is that the TM system integrates with the TMS.

You most certainly know that EBMT (Example-Based Machine Translation) seems to be one of the buzzwords in the industry. Some translation software already support something like "assembling", others are keen to do so. So, in the future, when working with TM technology, it will be an invaluable asset to have quality terminology as the TM also assembles from there where it finds no exact or fuzzy matches.

Mark D. Childress remacha:

I've heard this discussion at a conference or two recently. It seems to me all terminologists shudder at the thought but some have problems communicating why the sole use of translation memories is Not A Good Thing, whereas their use together with term databases is.

The simplest way to explain to those people not familiar with either, but who make the decisions whether to support terminology management or not:

- Translation memory: Descriptive — The way things are
- Terminology database: Prescriptive — The way things should be

Then ask them: Should one use a term that is (potentially) incorrect, simply because that's the way it is? Because if that's so, a term use error continues to multiply because there's no prescriptive ruling to use as a reference to standardize the term. As Kara mentioned in her post, it will ultimately make the translation memory useless.

Alternately, one can correct the problem each time it appears in "reality" and count the cost of correction, whereas defining the "ideal" goes a long way towards reducing the occurrence of the problem to begin with. Basic quality rule: Get it right the first time, and reduce your follow-on costs.

Propuesta

El proceso que ha iniciado Javier Fernández Serrador de unificación de las terminologías de los distintos proyectos de traducción de interfaces de usuario (Gnome, KDE, es@li, LuCAS) es fundamental. Propone la creación de una base de datos que incluya las traducciones consensuadas de los términos, y aquellas que se consideren incorrectas.

Cuando las herramientas estén maduras para utilizarlas (pienso en los desarrollos deseables de gtranslator y kbabel) implementar bases de datos terminológicos en TBX.

Correctores gramaticales y de estilo

¿Tengo que explicar que un corrector gramatical necesita una gramática?

Estudiar *diction*. Aplicarlo al español.

Corpora monolingües

Discusión

Según Catherine Ball (<http://www.georgetown.edu/cball/corpora/> [<http://www.georgetown.edu/cball/corpora/>]):

To create a rudimentary *concordancer* is a simple programming task: it is a matter of indexing words to lines, sorting the words alphabetically, and displaying each word in a fixed amount of context. However, most of the generally-available concordancers have many additional features, including options for producing full or partial concordances, sorting in a variety of orders, searching for collocations, and producing basic text statistics.

Propuesta

1. *Herramientas estadísticas*. Como traductor opino que necesitamos concordancias (cf. XConcord, de Cíbola; tea, de Masao Utiyama), colocaciones, análisis de patrones léxicos y agrupaciones de palabras (cf. Paai's text utilities)
2. *Anotaciones* (taggers). Mediante etiquetado XCES si pretendemos avanzar en el campo de la traducción automática. Desarrollar/liberar herramientas de desambiguación
3. Extracción de términos
4. Diferenciación entre cadenas que deben ser traducidas y las que no (XLIFF)

Corpora paralelos y multilingües

Discusión

Es evidente el interés de estas herramientas. Condiciones previas son

1. la *segmentación*, el *recuento de palabras* y
 2. el *alineado* de los textos paralelos, en una primera fase a nivel de párrafo (insuficiente). Comprender la distinción entre alineación oracional y léxica. Lucha contra paralelizaciones ruidosas y ambigüedad sintáctica y/o semántica.
- WA Gale y KW Church sientan las bases teóricas (estadísticas) sobre el alineado de texto en “A Program for Aligning Sentences in Bilingual Corpora”, de 1993 (<http://citeseer.nj.nec.com/gale93program.html> [<http://citeseer.nj.nec.com/gale93program.html>]; el documento incluye la fuente del programa)
 - XAlign, de Cíbola
 - Aportaciones del proyecto MULTTEXT (multext_align)
 - Desarrollos matemáticos a la última que escapan a mi comprensión actual, en las páginas de Dan Melamed (<http://www.cs.nyu.edu/~melamed/interests.html> [<http://www.cs.nyu.edu/~melamed/interests.html>])

Melamed ha creado también gran número de utilidades libres.

Propuesta

Segmentación y alineado a nivel de párrafo están favorecidos por el marcado DocBook de nuestros documentos.

El alineado a nivel de párrafo no es suficiente. Ver SRX, de Oscar/LISA.

Nuestro problema es definir las Unidades de Traducción.

Creación de memorias de traducción: «un corpus alineado y anotado constituye una memoria de traducción» (Abaitua, Tradumática 0, 2001).

Memorias de traducción

Discusión

Las TM contribuyen a la calidad de las traducciones en dos aspectos: su coherencia y su exactitud.

Propuesta

gtranslator y kbabel utilizan de forma natural MT cuando trabajan con ficheros .po.

Deben extenderse para utilizar MT cuando trabajen con documentación y textos libres. Para ello es necesario haber resuelto los problemas de la segmentación, alineación y marcado.

Los proyectos deben trabajar (o al menos ser capaces de importar y exportar) el formato TMX.