

Wikipedia eta itzulpen automatikoa: «harri batez bizpalau xori»

*Iñaki Alegria¹, Unai Cabezón¹, Unai Fernandez de Betoño²,
Gorka Labaka¹, Aingeru Mayor¹, Kepa Sarasola^{1*}, Arkaitz Zubiaga²*

¹ Ixa Taldea, <https://ixa.si.ehu.es>

² Euskal Wikipedia, <http://eu.wikipedia.org>

*kepa.sarasola@ehu.es

Jasoa: 2013-06-24

Onartua: 2013-10-14

Laburpena: Artikulu honetan elkarlanean egindako proiektu bat aurkezten dugu. Boluntario talde bat bildu dugu espainierazko Wikipediako hainbat artikulu euskarara itzultzeko, baina boluntarioen lana errazteko, Matxin itzultzaile automatikoa erabili dugu aurreitzulpenak sortzeko, eta horrela boluntarioen lana erre eta akatsak dituzten itzulpen automatiko horiek aztertu eta zuzentzea izan da. Lan honekin, batetik, Euskal Wikipedia aberastu dugu, 50.000 hitz berri gehituz. Beste alde batetik, sistema automatikoaren itzulpenak eta posteditatutako bertsio zuzenduekin corpus bat sortu dugu. Corpus hori erabili dugu posteditore estatistiko bat sortzeko, Matxin itzulpen automatikoko sistemaren irteeraren doitasuna % 10ean hobetuz.

Hitz gakoak: Wikipedia, itzulpen automatikoa, corpus.

Abstract: In this paper we define a collaboration framework that was tested with editors of Basque Wikipedia. Their post-editing of Computer Science articles has been used to improve the output of a Spanish to Basque MT system called Matxin. For the collaboration between editors and researchers, we selected a set of 100 articles from the Spanish Wikipedia. These articles would then be used as the source texts to be translated into Basque using the MT engine. A group of volunteers from Basque Wikipedia reviewed and corrected the raw MT translations. This collaboration ultimately produced two main benefits: (i) the change logs that would potentially help improve the MT engine by using an automated statistical post-editing system, and (ii) the growth of Basque Wikipedia. The results show that this process can improve the accuracy of an Rule Based Machine Translation system in nearly 10% benefiting from the post-edition of 50,000 words in the Computer Science domain.

Keywords: Wikipedia, automatic translate, corpus.

1. SARRERA

Artikulu honetan elkarlanaren onuraz arituko gara. Lan asko egin behar denean eta baliabideak urriak direnean, indarrak biltzea izaten da gakoa. Eta lan bera erabiliz hainbat behar asetzen baditugu edota hainbat helburu betetzen baditugu, askoz hobe. Gurean bi behar, bi helburu, izan ditugu eskuartean: itzulpen automatikoaren garapena [1] batetik eta Euskal Wikipediaren¹ aberastea bestetik.

Euskal Herriko Unibertsitateko IXA taldean² Hizkuntzaren Prozesamenduaren arloan ikerketa egiten dugu [2,3], besteren artean itzulpengintza automatikoan. Euskararekin lan egiten duen eta publikoki erabilgarria izan zen lehenengo itzulpen automatikoko sistema, Matxin, eraiki ondoren, erregeletan oinarritutako sistema horren irteera hobetzeko posteditore estatistiko bat eraikitzea erabaki genuen. Baina horretarako automatikoki itzulitako esaldiak eta horien eskuzko postedizioak bilduko lituzkeen corpus bat beharrezkoa genuen.

Lankidetzaz sortutako Wikipedia entziklopedia eleaniztuna euskaraz ere badugu, baita osasun osoz eduki ere [4,5]. Baina euskara bezalako baliabide urrietako hizkuntzetako Wikipediak aberastea ez da erraza, sarrera berriak idatziko dituzten boluntarioen kopurua txikia delako, eta beraz, boluntario bakoitzak egin beharreko esfortzua handia delako. Editore kopuru txikia duten hizkuntzek ezin dute lehiatu ingelesa edo espainiera bezalako hizkuntzetan izaten den Wikipediaren hazkuntzarekin. Baina alde positibotik ikusita, Wikipedia txikitako editoreak profitatu ahal izango dira hizkuntza handietan sortutako eduki kopuru handiez, eduki horiek haien hizkuntzara itzuliz, bide hau eduki berriak sortzea baino askoz ere merkeagoa izanik [6].

Itzulpen automatikoa laguntza ederra izan daiteke itzulpen-prozesu hori errazte aldera. Gaur egungo sistema automatikoez ematen dituzten itzulpenak okerrak eta akatsez betetakoak badira ere, gure hipotesiaren arabera itzultzaile ez profesionalen kasuan sistema automatikoaren irteera zuzentzea hutsetik itzultzea baino errazagoa da.

Hona hemen gure ideia: Euskal Wikipediako sarrera berriak idaztea espainierazko Wikipediako artikuluen itzulpen automatikoa zuzenduz. Horrela, alde batetik, editoreen lana erraztu dezakegu Euskal Wikipedia aberasteko eta, bestetik, itzulpen automatikoaren irteera okerrak eta eskuz zuzendutako postedizioak bildu ditzakegu corpus batean, eta postediziozko corpus horren bitartez posteditore estatistikoa eraiki gero [7,8]. Esaera zaharrak dioen bezala: «*harri batez bizpalau xori*».

¹ <http://eu.wikipedia.org>

² <http://ixa.si.ehu.es>

2. AURKEZPENA

2.1. Wikipedia: Entziklopedia askea

Wikipedia Interneten argitaratutako eduki askeko entziklopedia eleantza da. Lankidetzaz editatua, mundu osoko boluntarioek idazten dute. Gaur egun, Interneten dagoen kontsultarako tresnarik handiena eta zabalena da, zalantzarik gabe, eta baita Web 2.0 eta haren parte hartzearen filosofiaren eredu argienetakoa ere. Etengabe eguneratzen, handitzen eta zuzentzen da. Wikipediak ezagutza librea sustatzen du. *Libre* izatea ez da doakoa delako bakarrik, baita bildutako jakintzaren berrerabiltzea baimentzen duelako ere; jakintzaren eboluzioa sustatzeko oinarri sendoa izan da hori.

2001. urtean sortu zen eta hamar urtebete eta gutxira, 2012ko azaroaren lehenean, 285 hizkuntzako edizioak zituen, guztira 23 milioi artikulua, horietatik lau milioitik gora ingelesez. Milaka boluntario, profesional gutxi batzuk baino etekin hobea lortzen ari dira, kantitatea kalitate bihurtuz. Euskal Wikipedia, 150.000 artikulurekin, osasun onean dagoela esan genezake. Wikipediako hizkuntzen zerrendan euskara 35. posizioan agertzen da. Hori bai, tamalez, bere tamaina oraindik txikia da Hizkuntzaren Prozesamenduko aplikazio aurreratuetan corpus gisa erabilia izateko.



1. irudia. Euskal Wikipedia.

Esan beharra dago euskaraz sekula ez dela izan Wikipedia baino entziklopedia zabalagorik. Euskal Wikipediako komunitatea lanean gogoz dabil, bai euskarazko artikulua berriak sortzen, bai editore berriak bilatzen. Eta lan honetako helburu bat horixe izan da: Euskal Wikipedia aberastea.

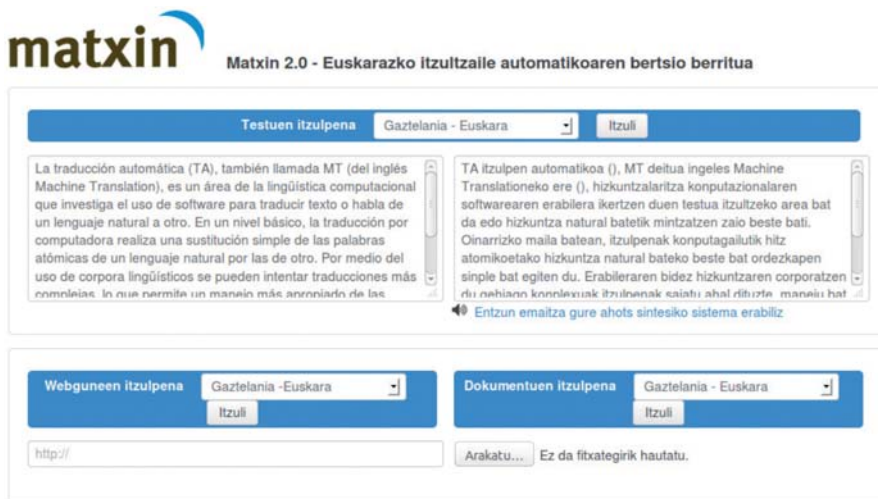
2.2. Matxin: Itzultzaile automatikoa

Matxin [9,10] erabilera publikoko euskararako lehenengo itzulpen automatikoko sistema izan zen, Euskal Herriko Unibertsitateko IXA Taldearen eta Elhuyarren artean garatua. Software libre da, GNU GPL lizentziapean argitaratutakoa.

Matxin erregeletan oinarritutako itzulpen automatikoko sistema bat da. Hiru fasetan egiten du itzulpena: analisia (espainierarako Freeling [11] software libreko paketea erabiltzen du), transferentzia (lexikala eta estrukturala) eta sorkuntza (sintaktikoa eta morfologikoa).

Ezberdintasun lexikal eta sintaktiko handiak aurkezten dituzten hizkuntzen arteko itzulpenak egiteko eta abiapuntu-hizkuntzatik independentea izateko diseinatuta dago. Matxin 2.0 prototipoak espainieratik euskarara itzultzen du eta erabilera orokorrekoa da. Sarean erabil daiteke bere webgunean³ eta OpenTrad proiektuko webgunean⁴ proba daiteke, eta Sourceforge-n kode irekiko software libre gisa banatzen da⁵.

Gaur egun espainiera-euskara sistema hobetzen jarraitzeaz gain, ingelesetik euskarara itzultzen duen prototipoa garatzen ari da IXA taldea. Gainera, estatistikan oinarritutako sistemetan eta sistema hibridoetan ere ikerketa egiten ari da [12].



The screenshot shows the Matxin 2.0 web interface. At the top left is the 'matxin' logo. To its right is the text 'Matxin 2.0 - Euskarazko itzultzaile automatikoaren bertsio berria'. Below this, there are several sections:

- A header bar with 'Testuen itzulpena', a dropdown menu set to 'Gaztelania - Euskara', and an 'Itzuli' button.
- A text area containing a description of automatic translation (TA) in Spanish: 'La traducción automática (TA), también llamada MT (del inglés Machine Translation), es un área de la lingüística computacional que investiga el uso de software para traducir texto o habla de un lenguaje natural a otro. En un nivel básico, la traducción por computadora realiza una sustitución simple de las palabras atómicas de un lenguaje natural por las de otro. Por medio del uso de corpora lingüísticos se pueden intentar traducciones más complejas. Lo que permite un manejo más avanzado de las...'
- Another text area containing a description of automatic translation in Basque: 'TA itzulpen automatikoa (), MT deliaua ingeles Machine Translationeko ere (), hizkuntzalaritza konputazionalaren softwarearen erabilera ikertzen duen testua itzultzeko area bat da edo hizkuntza natural batetik mintzatzen zaio beste bati. Oinarrituko maila batean, itzulpenak konputagailutik hitz atomikoetako hizkuntza natural bateko beste bat ordezkapen simple bat egiten du. Erabileraren bidez hizkuntzaren korporatzen...'
- A button labeled 'Entzun emaitza gure ahots sintesiko sistema erabiliz'.
- A section for 'Webguneen itzulpena' with a dropdown menu set to 'Gaztelania - Euskara' and an 'Itzuli' button.
- A section for 'Dokumentuen itzulpena' with a dropdown menu set to 'Gaztelania - Euskara' and an 'Itzuli' button.
- A text input field containing 'http/'.
- An 'Arakatu...' button and the text 'Ez da fitxategirik hautatu.'

2. irudia. Matxin itzulpen automatikoko sistema on-line.

³ <http://matxin.elhuyar.com>

⁴ <http://www.opentrad.com>

⁵ <http://matxin.sourceforge.net>

Itzulpen automatikoko sistemak hobetzeko modu bat postedizio-lanaz baliatzea da. Postedizioa sistema automatikoa ematen dituen itzulpenen zuzenketa da. Postedizioa gizakiek egin ohi dute, sistema automatikoaren kalitatea ona denean itzulpenak egiteko lana murriztu baitaiteke horrela.

Posteditore automatikoa ere eraiki daitezke. Horretarako itzultzaile automatikoaren emaitzak eta emaitza horien eskuz posteditatutako irteerak erabiltzen dira, horiekin ikasketa automatikoa aplikatu ahal izateko [13]. Hauxe da, hain zuzen, elkarlan hau abiarazi zuen asmoa: Matxin erregeletan oinarritutako sistemaren irteera hobetuko lukeen posteditore estatistiko bat eraikitzea.

2.3. **OmegaT: Ordenagailuz lagundutako itzulpen-tresna**

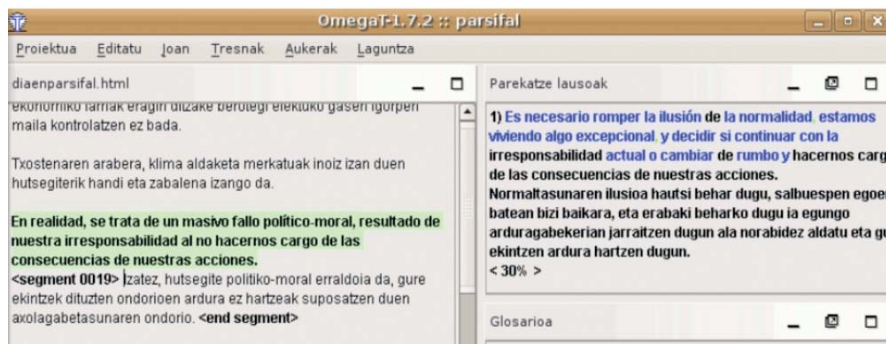
OmegaT⁶ aplikazioa ordenagailuz lagundutako itzulpen-tresna bat da, giza-itzultzaileei laguntzen diena, besteak beste, itzulpen-memoriekin. OmegaT-k, aurretik egindako hainbat itzulpen gordeta dituzenez, itzultzaileak esaldi berri bat itzuli behar duenean, aurretik itzulitakoen artean antzekoenak bilatu eta haien artean zeuzkan gordetako itzulpenak proposatzen ditu, baten bat egokia izanez gero giza-itzultzaileak berrerabil dezan. Itzulpen-memoriak oso erabilgarriak izan daitezke itzuli beharreko testuen artean antzeko atalak dituzten testu-zatiak agertzen badira.

Itzulpen-memoriaz gain, OmegaT-k beste hainbat funtzionalitate eskaintzen ditu:

- Internetetik zuzenean hainbat itzultzaile automatiko erabiltzeko aukera eskaintzen du (adibidez, Google Translate, Apertium edo Belazar), jatorrizko esaldien oinarritzko itzulpen bat eskuratzeko.
- Glosario eta hiztegien sorkuntza eta inportazioa. Terminologiaren kudeaketarako laguntza da hori, itzulpenen koherentzia lortzeko. Glosarioetan domeinu berezi baterako hitz bakanak edo hitz anitzeko unitate lexikalak gordetzen dira. OmegaT-k itzuli behar dugun uneko segmentuko hitzen ordainak erakusten ditu, hiztegiaren edo glosarioaren baldin badaude.
- Lematizatzaileen erabilera, hitzen lema hobeto eratzeko. Modu horretan asko hobetzen da parekatze lausoaren bilaketa eta glosarioen erabilera. Ezaugarri hau euskara bezalako hizkuntza deklinatuen zati bereziki interesgarria da.

OmegaT-k hainbat formatutako fitxategiak onartzen ditu (HTML, MS Office eta OpenOffice bulegotika-aplikazioen formatuak, DocBook, PO, etab.) eta fitxategiak iragazteko dituen arauak jarraituta formatu-markak identifikatu eta itzuli behar den testu soila lortzen du. Testu gordin hori

⁶ <http://www.omegat.org>



3. irudia. OmegaT tresnaren interfazea.

segmentutan banatu eta segmentuak banaka erakusten dizkio giza itzultzaileari beronek lehen aipatutako laguntzekin itzul ditzan.

OmegaT kode irekiko plataforma da eta erabiltzaileen eta garatzaileen komunitate aktibo batek sostengatzen du.

OmegaT itzulpen-ingurunea aukeratu dugu Wikipediako espainierazko artikuluen itzuFWLlpen-automatikoen posteditzio lanetarako, software libre izatean, aukera emango digulako sistema gure beharretara egokitzeko, besteak beste funtzionalitate berriak integratuz.

Proiektuaren hasieran itzulpenean laguntzeko beste sistema posible batzuk aztertu genituen: (1) World Wide Lexicon Translator (WWL3)⁷, Firefox-erako gehigarri bat webguneak itzulita ikusi ahal izatekoak, konbinatzen zuen giza-itzulpena eta itzulpen automatikoa, baina posteditatzeko interfazea ez zebilen oso ondo; (2) Google Translation Toolkit⁸ tresnak Wikipediako sarrerak itzultzeko laguntza eskaintzen du baina software libre eta irekia ez denez ezin izan genuen moldatu gure beharretara. OmegaT aukeratu genuen, software libre izanda moldatzeko aukera ematen zigan eta.

3. ELKARLAN PROIEKTUA

3.1. Diseinua

Gure proiektua boluntarioen lankidetzan oinarritu da. OmegaT plataforma egokitua erabiliz boluntarioek espainierazko Wikipediako artikuluen itzulpen-zirriborroa lortzen zuten Matxin itzultzaile automatikoa erabiliz

⁷ <http://wordpress.org/plugins/speaklike-worldwide-lexicon-translator>

⁸ <http://translate.google.com/toolkit>

eta, ondoren, itzulpen gordin hori aztertu eta zuzentzen zuten, Postedizio-lan horretan, batzuetan, itzulpen automatikoaren emaitza oso txarra zenean, itzulpen hori ez zen erabiltzen eta esaldi osoa berriro itzultzen zuen boluntarioak bere kabuz. Baina, gehienetan, automatikoki sortutako itzulpeneko akatsak zuzenduz lan asko aurrezten zen.

Matxin itzulpen automatikoko sistemaren emaitzak hobeak izateko asmoz ezagutza-eremu zehatz baterako egokitzea erabaki genuen. Aukeratu genuen eremua informatikarena izan zen, batetik eremu teknikoa izatean faktore kulturalen menpekotasuna txikiagoa delako eta bestetik gure gertuko komunitatearentzat gai ezaguna delako. Matxin sistemaren lexikoiak eremu horretarako egokitu genituen eta, noski, *Informatika* kategoriakoak izan ziren gure proiektuan itzultzeko aukeratuko genituen Wikipediako artikulak.

Proiektua publiko egin aurretik eta boluntarioak biltzen hasi aurretik, probako saiakera batzuk egin genituen diseinatzen ari ginen prozesua aztertzeko, eta artikulua luzean postedizioa egitea esfortzu handiko lana zela ikusi genuen (artikulu luzeak posteditatzeko 8 ordu baino gehiago behar izaten ziren). Horrek boluntarioak lortzea zaildu zezakeenez eta lana bukatu gabe uzteko posibilitatea handitzen zuenez, artikulua luzeak itzuli beharrean neurri ertaineko artikulak proposatzea erabaki genuen.

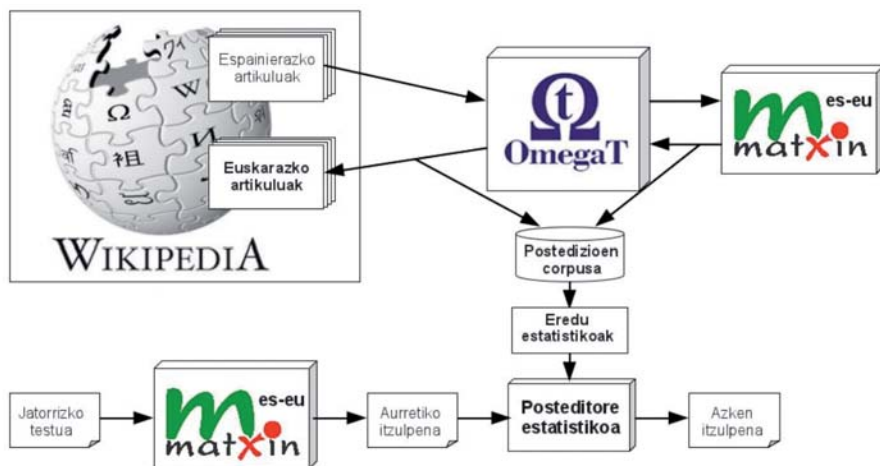
Boluntario bakoitzari hasieran 2-3 lerroko artikulua motz batekin probatzen hastea eskatu zitzaion, OmegaT plataforma erabiltzeko prozesu osoan trebatzeko: artikulua Wikipediatik jaistea, itzulpen automatikoa lortzea, itzulpena posteditatzeko eta azkenik emaitza Euskal Wikipediara igotzea. Ondoren 10-20 lerroko artikulua bat itzultzeko eskatzen zitzaion, eta gustura aritu ziren batzuk artikulua bat baino gehiago ere itzuli zuten.

Wikiproiektu bat martxan jarri genuen⁹. Bertan OmegaT-ren bertsio egokituja jaisteko loturarekin batera, Espainerazko Wikipediako hainbat artikuluren zerrenda jarritz, garai horretan euskarazko itzulpenik ez zutenak, boluntario bakoitzak artikulua horien artean gogokoena(k) aukeratu ahal izateko.

Lankidetzak-kanpaina publikoa zortzi hilabetez egon zen martxan, 2011ko uztailetik 2012ko otsailera. Boluntarioekin aritzeko trebatze-saio pare bat antolatu genituen, artikulua motzekin lortu behar zen trebakuntza hori taldean eta lagunduta landu ahal izateko. Gainera, «Wikitzul» deituriko elkarlan-saioetara boluntario gehiago erakartzeko bi deialdi egin ziren¹⁰.

⁹ http://eu.wikipedia.org/wiki/Wikiproiektu:OpenMT2_eta_Euskal_Wikipedia

¹⁰ <http://www.unibertsitatea.net/blogak/ixa/2011/12/05/euskal-wikipediaren-edizio-maratoia-durangoko-azokan>



4.irudia. Prozesu osoaren ikuspegia.

Guztira, prozesu osoan, 36 boluntariok parte hartu zuten eta, 100 artikulua itzuli ondoren, euskarazko 50.204 hitz gehitu dira Euskal Wikipediara.

Eskuz posteditatutako artikulua horiekin bi helburu lortu ditugu. Batetik, Euskal Wikipedia zabaltzea eta, bestetik, automatikoki itzulitako testuen eta haien eskuzko postedizioen corpusa osatzea, posteditore estatistikoa bat eraikitzeke erabili dena. Gainera, prozesu honetan, bai Matxin sisteman bai OmegaT-n hainbat hobekuntza egin dira.

3.2. Wikipediako artikuluen aukeraketa

Boluntarioak gureganatu ahal izateko, trebakuntzarako artikulua labur batzuk, baina batez ere luzera ertaineko artikulua asko identifikatu behar genituen. Azkenean luzera ertaineko artikuluen itzulpenen postedizioa proposatzea erabaki genuen, hasierako kontaktuetan argi ikusi baikenuen proposatutako artikulua handiegia izanez gero ez zela boluntariorik hurbilduko.

Wikipediako kategoria batean dauden artikuluen zerrenda ematen duen programa prestatu genuen informatika eremuko artikulua ertainak eskuratzeko. Artikulu bakoitzean ea beste hizkuntzetan baliokideak ba ote zeuden aztertu genuen eta horien luzerak ere eskuratu genituen programa horrekin.

Katalanezko Wikipediaren tamaina (378.408 artikulua) espainierazko (902.113 artikulua) eta euskarazko (135.273 artikulua) Wikipediaren tamainaren artean kokatzen zenez, gure ustez katalanezko Wikipedian dagoen eta euskarazkoan ez dagoen artikulua bat lehenago gehitu beharko litzateke Euskal Wikipedian, katalanezko Wikipedian ez dagoen artikulua bat baino.



5. irudia. Wikitzul saioa (2012/10/12).

Beraz, Euskal Wikipediaren hutsuneak identifikatzeko orduan gure programaren emaitza erabili genuen irizpide honekin: Katalanezko Wikipedian *Informàtica* kategorian eta Espainierazko Wikipedian izanik, euskarazkoan existitu gabe espainierazkoan 10-20 lerro arteko luzera zuten artikuluenak identifikatzeko.

Horrela 140 artikulua identifikatu genituen, Euskal Wikipedian sartzeko artikuluen proposamen-zerrendan boluntarioei eskaintzeko.

3.3. Matxinen egokitzapenak

Informatika eremurako egokitu dugu Matxin itzulpen-sistema, bere lexikoi elebiduna bi modutara aberastuta [14]:

- Lexikoiaren egokitzapena hiztegi elektronikoak erabiliz. Hainbat espainiera-euskara on-line hiztegitan informatika eremurako hitzen ordainen bilaketa sistematikoa egin dugu. Horrela 1.623 sarrera berri gehitu ziren Matxinen lexikoian. Termino gehienak hitz anitzekoak ziren, adibidez «base de datos» (datu-base) edo «lenguaje de programación» (programazio-lengoaia). Hitz bakarreko termino batzuk ere lortu ziren, adibidez «iterativo» (iteratibo), «ejecutable» (exekutagarri) edo «ensamblador» (mihizatzaile). Gainera, hautapen lexikalerako ordainen ordena aldatu zen 184 hitzetan; esate baterako, erdarazko «rutina» sarreraren ordainen ordena aldatu genuen, «errutina» ordaina jarri genuen lehenengo ordain modura, ordura arte lehenesten zen «ohitura» ordaina atzerago eramanda.

- Lexikoia-aren egokitzapena corpus paralelo bat erabiliz. Informatika-aren eremuko espainiera-euskara corpus paralelo bat bildu genuen, Mozilla software librearen lokalizazioan sortutakoa (138.000 segmentu, 600.000 hitz espainieraz eta 440.000 euskaraz). Corpus hau itzulpen automatiko estatistikorako egokia ez bada ere, erabilgarria izan daiteke erlazio lexikalak erauzteko. Giza++ lerrotzetan oinarrituta, corpuseko sarrera bakoitzerako itzulpen posibleen zerrenda erauzi genuen, itzulpen posible bakoitzaren probabilitatearekin. Doitasunagatik, zerrenda hauek hautapen lexikalerako soilik erabili genituen. Ordainen ordena aldatu zen lexikoiko 444 sarreratan. Adibidez »dirección» sarrerarako «helbide» hobetsi zen, «nora-bide» beharrean, egokiagoa izango zelakoan.

3.4. OmegaT-ren egokitzapenak

Posteditoreentzat erabilerrazago egiteko, OmegaT egokitu dugu hainbat funtzionalitate gehituz:

- Matxin es-eu itzultzaile automatikoaren integrazioa. OmegaT-k badu itzulpen automatikoko zerbitzuak konektatzeko klase bat, zerbitzu berriak erraz gehitu ahal izateko. Matxin OmegaT-ren barruan integratzeko erabili genuen guk klase hori. Integrazio-lana errazteko, Matxin sistema egokitu genuen web zerbitzu gisa inplementatuz, SOAP bidez API deiekin atzigarria izateko. Horrela OmegaT barrutik erraz erabil daiteke Matxin sistema.
- Inportatu/esportatu Wikipediako artikulua OmegaT-n. OmegaT-k eskaintzen dituen MediaWiki dokumentuen inportaziorako aukerei gehituz, ezaugarri berri bat inplementatu genuen espainierazko Wikipediako artikulua inportatzeko eta beste bat posteditatutako artikulua Euskal Wikipediara igotzeko. Horretarako *login* egiteko modulu berri bat inplementatu behar izan genuen. Artikulu berria Euskal Wikipediara igotzean, gainera, Matxinek emandako itzulpen automatikoa eta boluntarioak posteditatutako itzulpena geure zerbitzari batera bidaltzen dira, lortu nahi dugun postedizio-corpusa osatzen joateko.
- Euskarazko zuzentzaile ortografikoaren integrazioa, postedizio-lanak errazteko.
- Wikipediako estekak itzultzeko programa, Wikipediako metadatuaren informazioa erabiliz. Ikus dezagun, adibidez, `[[gravedad|gravedad]]` Wikipediako barne loturaren itzulpena. Lotura horren lehenengo terminoak Wikipediako sarrera bati egiten dio erreferentzia, eta bigarren terminoa lotura horretan erakutsiko den testua da. Wikipediako artikulua bati erreferentzia egiten dion lehenengo termino hori itzultzeko, gure programak Wikipedia barruko informazioa erabiltzen

du (zein beste hizkuntzatan dagoen artikulua eta zein den bere sarrera hizkuntza horietan) espainierazko artikulua horri dagokion euskarazko artikulua eskuratzeko, kasu honetan «grabitazio». Bigarren terminoa, lotura hori adierazteko erakutsiko den testua, Matxin itzultzaile automatikoa erabiliz itzuliko da, «larritasuna» lortuz. Horrela, emandako barne loturaren itzulpena [[grabitazio | larritasuna]] izango da. Posteditorearentzat laguntza ederra da automatikoki lortzea Euskal Wikipediako lotura (*grabitazio*), berak ezer bilatu beharrik gabe. Gainera, euskaraz dagokion sarrera ikusita erraz zuzendu dezake Matxin-ek aukeratu duen ordaina, testuinguru horretan egokia ez bada (*larritasuna-->grabitazioa*). Ordezkapen hori beti automatikoki egitea lagungarria izan daiteke, baina ez beti; epe erdirako eginkizunen artean jarri dugu funtzionalitate hau aztertu eta inplementatzea.

4. EMAITZAK ETA HOBEKUNTZAK

Guztira hauek dira sortu eta modu irekian plazaratu ditugun produktuak:

— Sortutako Corpusak:

- Espainiera/euskara corpus paralelo bat.¹¹ Mozilla softwarearen lokalizazioan sortu denaren bertsio berria (138.000 segmentu, 600.000 hitz espainieraz eta 440.000 euskaraz).
- Testu itzuli eta zuzenduen corpus bat.¹² Espainierazko Wikipediako 100 artikulua Matxin itzultzaile automatikoak sortutako itzulpenekin eta boluntarioek eskuz posteditatutako itzulpen zuzenduekin. Corpus horren euskarazko aldeak 50.204 hitz dauzka.

— Wikipedia:

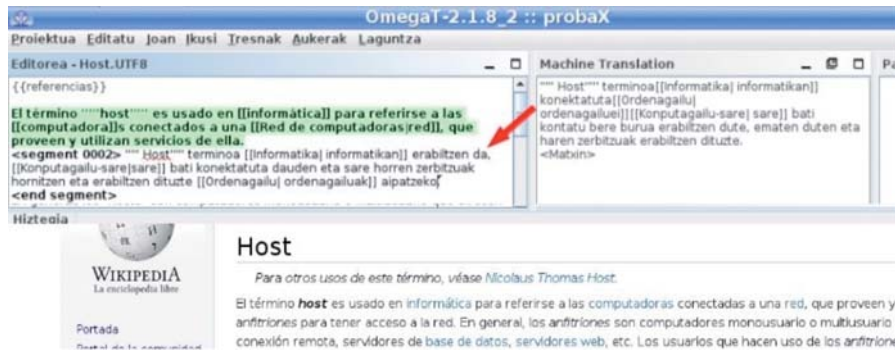
- Euskarazko artikulua berriak: 100 artikulua berri gehitu dira Euskal Wikipedian, guztira 50.204 hitz.¹³
- Artikuluak bilatzeko programatxoa (wikigaiak4koa.pl). Tresna hau, perl-ez inplementatua, Wikipediaren edozein hizkuntzatarako kategoria baten edukia aztertzeke erabil daiteke. Kategoria bat eta lau hizkuntza emanda, lehenengo hizkuntzaren kategoria horretako artikuluen zerrenda ematen digu, artikulua bakoitzaren beste hiru hizkuntza horietako baliokideekin eta beren luzerekin.¹⁴

¹¹ <http://ixa2.si.ehu.es/glabaka/lokalizazioa.tmx>Eskerrak Elhuyarri eta Julen Ruiz-i

¹² <http://ixa2.si.ehu.es/glabaka/OmegaT/OpenMT-OmegaT-CS-TM.zip>

¹³ <http://eu.wikipedia.org/w/index.php?title=Berezi:ZerkLotzenDuHona/Txantilo:OpenMT-2&limit=250>

¹⁴ <http://www.unibertsitatea.net/blogak/testuak-lantzen/2011/11/22/wikigaiak4koa>



6. irudia. Wikipediako esteken itzulpena. Espainerazko Wikipediako «Host» artikuluan «red» hitza esteka moduan agertzen da. Esteka hori [[Red de computadoras | red]] moduan adierazten da. OmegaT eta Matxinek [[Konputagailu-sare|sare]] eskaintzen dute.

— Matxin:

- Matxinen bertsio egokitua informatika arlorako. Sistemaren lexikoa eremu zehatz horretara egokitu dugu eta SOAP¹⁵ zerbitzu moduan implementatu dugu. Jatorrizko Matxin sistema eta informatika arlorako egokitutakoa automatikoki ebaluatu ditugu. Espainierazko Wikipediatik aukeratutako artikulua erabili dira, boluntarioek postedizio bidez sortutako itzulpen zuzenduak erreferentzia gisa hartuta. Erabilitako metrika automatiko [15] guztietarako (MBLEU, BLEU, NIST, METEOR, TER, WER eta PER) egokitutako sistemak emaitza hobekuntza ematen ditu: hobekuntza erlatibo handienak BLEUk eta MBLEUk adierazten dute (% 15) eta hobekuntza txikiak WER metrikak (% 3,5)
- Postedizio automatikorako modulua. Automatikoki lortutako itzulpenak eta itzulpen horien eskuz posteditatutakoak biltzen dituen 50.000 hitzeko corpusa baliatuta eta teknika estatistikoak erabiliz Matxin sistemarekin lortutako itzulpenak automatikoki posteditatzen dituen programa berri bat sortu dugu. Programa honek Matxin sistemaren irteera jaso eta postedizio automatikoa egiten du. Ebaluazio automatikoak erakutsi du sistema berri honek % 10eko hobekuntza lortzen duela, Matxin sistema soilarekin konparatuz gero [14]. Gurean erabilitako corpusaren tamaina postedizio estatistikoari buruzko nazioarteko esperimendu nagusietan erabili dituztenak baino txikiagoa da (adibidez Simard *et*

¹⁵ <http://eu.wikipedia.org/wiki/SOAP>

al. [13] 100.000 hitzeko corpora erabiltzen dute). Beraz postedizio corpus handiago bat lortuz hobetzeko bidea badagoela aurreikusten dugu.

— OmegaT:

- Matxin itzultzailea eta euskarazko zuzentzailea OmegaT-n erabiltzeko aukera gehitu dugu.
- Wikipediako artikuluak inportatzeko eta esportatzeko funtzionalitateak gehitu ditugu. Ezaugarri hau hizkuntzarekiko independentea da, eta euskara ez beste hizkuntzentzat ere erabil daiteke (ezaugarri hau ez da oraindik probatu beste karaktere multzo bat erabiltzen duten hizkuntzetarako, arabierarako adibidez).
- Wikipediako estekak itzultzeko programa inplementatu dugu, Wikipedia barruko informazioa erabiliz espainierazko artikulua bati dagokion euskarazkoa zein den lortzen duena.
- OmegaT deskargatzeko, instalatzeko eta erabiltzeko eskuliburua sortu dugu,¹⁶ Wikipediako artikuluak posteditatzeko zehaztasunekin.

5. ONDORIOAK ETA ETORKIZUNEN LANA

Boluntario-lana lortzea zaila izan da gurean, egin beharreko lana egiteko komunitatea sortzea eta koordinatzea esfortzu handia suposatu baitugu. Bagenekien hasieratik euskara bezalako hiztun gutxi hizkuntza baten kasuan hala izango zela, eta hala ere gure helburua lortu dugu. 36 boluntariok parte hartu dute proiektuan eta horietako 20k luzera ertaineko artikulua bana amaitu dute. Argi gelditu zaigu Wikipediako artikulua motzak aukeratzea egokiagoa dela boluntarioak erakartzeko proiektura, gehienetan motibazioa ez delako nahiko handia gehiegizko esfortzua inbertitzeko.

Itzulpen automatikoak itzultzaile profesionalentzat oso lagungarria ez bazirudien ere, itzultzaile amateurrentzat erabilgarria izan zitekeela zen gure hasierako hipotesia. Proiektuan sortutako postedizioak aztertu ondoren, gure hipotesia baieztatu dezakegu, itzulpen automatikoko sistemaren irteeraren kalitatea oso ona ez bada ere, nahikoa delako editoreei laguntzeko, egin beharreko esfortzua gutxituta.

OmegaT lokalki instalatu eta konfiguratu behar izatea eragozpen bat izan da, posteditoreen komunitate handi bat erabiltzea nahi genduela kon-

¹⁶ http://siuc01.si.ehu.es/~jipsagak/OpenMT_Wiki/Eskuliburua_Euwikipedia+OmegaT+Matxin.pdf

tuan harturik. Gure proiekturako eskertzekoa izango zen Google Translation Toolkit bezalako on-line elkarlan plataforma bat erabili ahal izatea. Hau dela eta, etorkizuneko proiektuetan erabiltzeko on-line plataforma egokiago bat egokitzea edo sortzea planeatzen ari gara.

Erabiltzaile berri batek OmegaT tresna erabiltzeko izan ditzakeen zailtasunak ere kontuan hartu behar dira, plataformak aukera eta funtzionalitate asko eskaintzen dituelako, eta horien artean norberak behar duen azpimultzoa soilik erabiltzea konplexua egiten delako. Zorionez, dokumentazio asko existitzen da horretan laguntzeko; gainera, guk geure proiektuan erabiltzen diren funtzionalitateak argitzen dituen eskuliburu bat sortu dugu posteditoreen lana gidatzeko. Gure proiektuan ikusi ahal izan dugu, gida hauekin, erabiltzaileek hasierako zailtasunak gainditzen dituztela eta berehala lortzen dutela OmegaT-rekin bere kabuz lan egiteko trebetasuna.

Wikipediako artikuluen metadatuaren tratamendua erronka bat da, bai itzulpen automatikoko sistemarentzat eta bai giza itzultzaileentzat. Honetan lagungarria izan da guk inplementatutako programa, artikuluetan agertzen diren loturen baliokidetzak lortzen dituen, Wikipediaren hizkuntzar-teko loturen informazioa erabiliz. Etorkizunean Wikipediako esteketan eta barne antolakuntzan dagoen informazioa era sakonagoan erabiltzea aurreikusten dugu.

Buruan dugu itzulpen-sistemaren lexikoa aberastea ere, domeinuaren arabera ordain egokiagoak hautatzeko. Informatika ez den beste arlo batera ere zabal genezake gure sistema.

Etorkizuneko proiektuetan boluntarioak erakartzeko estrategiak ere findu beharko genituzke, boluntario-nitxo berriak identifikatuz edo: unibertsitateak, hizkuntz-eskolak, euskaltegiak...

6. ESKER ONAK

Ikerketa hau bi erakundek finantzatu dute: Eusko Jaurlaritzak (Berba-tek proiektua, Etorrek deialdiko IE09-262) eta Espainiako Hezkuntza eta Zientzia Ministerioak (OpenMT2 proiektua, TIN2009-14675-C03-01). Elhuyarrek eta Julen Ruizek lagundu ziguten baliabideak biltzen Matxin itzultzailea informatika arlora egokitzen. Eta gure eskerrik beroena 36 boluntario kolaboratzaileei, eurak izan baitira proiektuaren emaitza arrakastua lortzea ahalbidetu dutenaki.

7. BIBLIOGRAFIA

- [1] HUTCHINS W.J. eta SOMERS H. 1992. *An introduction to machine translation*. Academic Press, London. <http://goo.gl/U0IbV>
- [2] ADURIZ I., ALEGRIA I., ARTOLA X., DÍAZ DE ILARRAZA A. eta SARASOLA K. 2011. «Teknologia garatzeko estrategiak baliabide urriko hizkuntzetarako: euskararen eta Ixa taldearen adibidea». *Linguamatica*, **3**, 13-31
- [3] HERNÁEZ I., NAVAS E., ODRIOZOLA I., SARASOLA K., DIAZ DE ILARRAZA A., LETURIA I., DIAZ DE LEZANA A., OIHARTZABAL B. eta SALABERRIA J. 2012. «The Basque language in the digital age/Euskara aro digitalean». *METANET White Paper Series*. Springer.
- [4] FERNANDEZ DE BETOÑO U. 2011. «Hamar urte jakintza libreba zabal-tzen». *Gaur8.info*. <http://goo.gl/TO1to>. 2011-01-14.
- [5] MUJIKAA. 2011. «Wikipedia: milioika artikulua ez ezik, zaleak eta kritikoak pilatu dituen kaxa erraldoia». *Gaur8.info*. <http://goo.gl/TO1to>
- [6] WAY A. 2010. «Machine translation». *The Handbook of Computational Linguistics and Natural Language Processing*. Wiley-Blackwell, Oxford. 531-573
- [7] ALEGRIA I., CABEZÓN U., FERNANDEZ DE BETOÑO U., GONZALEZ G., ITURBE M., LABAKA G., MAYOR A., SARASOLA K. eta ZUBIAGA A. 2013. «OpenMT2 eta Euskal Wikipedia wikiproiektuaren emaitzak». *Proceedings of the IX. Informatika Euskaldunen Bilkura, IEB2013, Udako Euskal Unibertsitatea, Donostia*.
- [8] ALEGRIA I., CABEZÓN U., FERNANDEZ DE BETOÑO U., LABAKA G., MAYOR A., SARASOLA K. eta ZUBIAGA A. 2013. «Reciprocal Enrichment between Basque Wikipedia and Machine Translators». *The People's Web Meets NLP: Collaboratively Constructed Language Resources*. Springer.
- [9] MAYOR A. 2007. *Matxin: Erregeletan oinarritutako itzulpen automatikoko sistema baten eraikuntza estaldura handiko baliabide linguistikoak berrera-biliz*. Tesi-lana. LSI Saila (EHU). Donostia. <http://goo.gl/j5aBI>
- [10] MAYOR A., DIAZ DE ILARRAZA A., LABAKA G., LERSUNDI M. eta SARASOLA K. 2011. «Matxin, an open-source rule-based machine translation system for Basque». *Machine Translation Journal*, **25**, 1, 53-82.
- [11] ATSERIAS J., CASAS B., COMELLES E., GONZÁLEZ M., PADRÓ L. eta PADRÓ M. 2006. «FreeLing 1.3: Syntactic and semantic services in an open-source NLP library». *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC'06)*, 48-55.
- [12] LABAKA G. 2010. *EUSMT: Incorporating Linguistic Information into SMT for a Morphologically Rich Language. Its use in SMT-RBMT-EBMT hybridation*. Tesi-lana. LSI Saila (EHU). Donostia. <http://goo.gl/b0Rga>
- [13] SIMARD M., UEFFING N., ISABELLE P. eta KUHN R. 2007. «Rule-based translation with statistical phrase-based post-editing». *Proceedings of the Second Workshop on Statistical Machine Translation*, 203-206.

- [14] ALEGRIA I., DIAZ DE ILARRAZA A., LABAKA G., LERSUNDI M., MAYOR A. eta SARASOLA K. 2011. «Matxin-Informatika: Versión del traductor Matxin adaptada al dominio de la informática». *Proceedings of the XXVII Congreso SEPLN*, 321-322.
- [15] SNOVER M., DORR B., SCHWARTZ R., MICCIULLA L. eta MAKHOUL J. 2007. «A study of translation edit rate with targeted human annotation». *Proceedings of the 7th Biennial Conference of the Association for Machine Translation in the Americas (AMTA)*, 223-231.