

“OpenMT2 eta Euskal Wikipedia” wikiproiektuaren emaitzak

OmegaT itzulpen-tresna hobetuta
euskaraz eta Wikipediarekin aritu ahal izateko.



Euskal Herriko Unibertsitatea

IXA TALDEA

HIZKUNTZAREN PROZESAMENDUA



**Iñaki Alegria, Unai Cabezón,
Gorka Labaka, Aingeru Mayor, Kepa Sarasola**

Ixa Taldea <https://ixa.si.ehu.es>



WIKIPEDIA
Entziklopedia askea

**Unai Fernandez de Betoño, Galder Gonzalez,
Mikel Iturbe, Arkaitz Zubiaga**

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



Helburua

“OpenMT2 eta Euskal Wikipedia” wikiproiektua

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

Hiru elementu hauek integratzea:

- **OmegaT** itzulpenak egiteko ingurunea
- **Matxin**: Itzulpen automatikoa
- **Euskal Wikipedia**

Bide batez, hiru elementuotan
hobekuntzak egin nahi genituen.



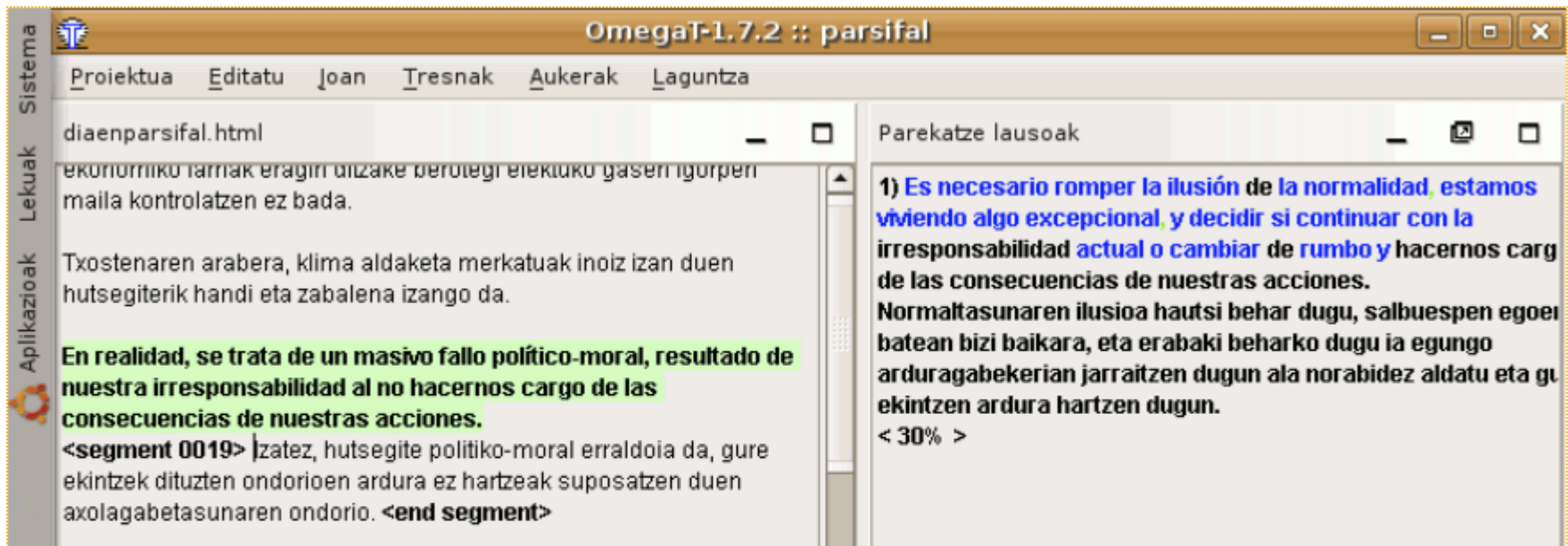
Hiru hankak

- **OmegaT** itzulpenak egiteko ingurunea
- **Matxin**: Itzulpen automatikoa
- **Euskal Wikipedia**

OmegaT



- **OmegaT** Ordenagailuz Lagundutako Itzulpen-sistema plataforma anitza da, kode irekikoa, parekatze lausoz, itzulpen-memoriaz, hitz gakoan bilaketaz eta glosarioz hornitutakoa, eta egindako itzulpenak proiektu eguneratuetan aprobetxatzea ahalbidetzen duena.





Matxin



- Hiru sistema automatiko erabil daitezke espainieratik euskarara itzultzeko:
 - Matxin
<http://www.opentrad.com/>
 - Itzultzailea (Eusko Jaurlaritza)
<http://www.itzultzailea.euskadi.net/traductor/portalExterno/text.do>
 - Google Translate
<http://translate.google.com/#es/eu/>
- Matxin
 - Lehenengoa (2006tik erabilgarri)
 - Software librea

Wikipedia



- Eduki askeko **entziklopedia**, lankidetzaz editatua, eleanitza, Interneten argitaratua
- Barne- eta kanpo-**esteka ugari** ditu.
- Programa informatikoen “**hodeian**” **erabil dezakete**
- Guztira 23 milioi artikulu, **hainbat hizkuntzatan**:
 - Ingelesezkoa: 4.000.000 artikulu
 - Espaineraz: 1.000.000
 - Katalanez: 400.000
 - **Euskaraz: 150.000** (35. hizkuntza; txikia HPko aplikazioetarako)
 - Galegoz: 100.000
 - ...



Hobekuntzak

- **OmegaT-n**
- **Matxin-en**
- **Euskal Wikipedia-n**

Hobekuntzak OmegaT-n (I)

- Matxin itzultzailea eta Xuxen erabili ahal izatea.

The screenshot shows the OmegaT-2.1.8_2 interface with the following elements:

- Editor:** Editore - chienWu.txt. The main text area contains a Spanish sentence: "Falleció debido a un segundo infarto, el [[16 de febrero]] de [[1997]]." Below it, there is a segment of text: `<segment 0204> Enlaces externos <end segment>` followed by a list of URLs and their descriptions.
- Parekatze lausoak (Translation Memory):** A window showing a match with the text: `1) == Enlaces externos ==`, `==Kanpo loturak==`, and `<100/100/33% lehenetsitakoak.tmx >`. A red circle highlights this window, and a red arrow points from the `Enlaces externos` text in the editor to it.
- Machine Translation:** A window showing the text: `Kanpoko estekak` and `<Matxin>`. A red circle highlights this window, and a red arrow points from the `<Matxin>` text in the editor to it.
- Glosarioa (Glossary):** An empty window.
- Hiztegia (Dictionary):** A window at the bottom left.
- Status Bar:** Shows the time "10:47 AM-ean automatikoki gordetako proiektua" and progress indicators "3/11 (18/281, 284)" and "16/16".

Hobekuntzak OmegaT-n (II)

Wikipediako estekak itzultzeko laguntza

The screenshot shows the OmegaT-2.1.8_2 software interface. The main window is titled "OmegaT-2.1.8_2 :: probaX". The interface is divided into several panes:

- Editorea - Host.UTF8:** Contains the source text in Spanish: "El término 'host' es usado en [[informática]] para referirse a las [[computadora]]s conectados a una [[Red de computadoras|red]], que proveen y utilizan servicios de ella." Below this is a translation snippet: "<segment 0002> 'Host' terminoa [[Informatika|informatikan]] erabiltzen da, [[Konputagailu-sare|sare]] bati konektatuta dauden eta sare horren zerbitzuak hornitzen eta erabiltzen dituzte [[Ordenagailu|ordenagailuak]] aipatzeko." A red arrow points to the word "Host" in the source text.
- Machine Translation:** Shows the translated text: "'Host' terminoa [[Informatika|informatikan]] konektatuta [[Ordenagailu|ordenagailuei]] [[Konputagailu-sare|sare]] bati kontatu bere burua erabiltzen dute, ematen duten eta haren zerbitzuak erabiltzen dituzte." Below this is a "Glosarioa" (Glossary) section.
- Wikipedia Article:** The bottom part of the image shows the Wikipedia article for "Host". The title "Host" is highlighted. A yellow box highlights the first sentence: "Este artículo o sección necesita **referencias** que aparezcan en una **pub** páginas de Internet **fidedignas**." A red box highlights the phrase "a las computadoras conectadas a una red, que" and a blue box highlights "a una red, que".

Hobekuntzak OmegaT-n (III)

- Matxin itzultzailea eta Xuxen erabili ahal izatea.
- Wikipediako estekak itzultzeko laguntza:
 - [[red de computadoras|red]] ---> [[konputagailu-sare|sare]]
 - [[gravedad|gravedad]] ---> [[grabitazio|larritasuna]]
larritasuna edo grabitazioa??
- Wikipediako artikuluak inportatzeko eta esportatzeko funtzionalitate berriak

OmegaT programaren bertsio hobetua:

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



Hobekuntzak Euskal Wikipedian

Euskal Wikipedian sortu diren 100 artikulu berriak.

(50.000 hitz)

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

wikigaiak4koa.pl perl programa.

Wikipediako kategoria bateko artikuluaren lista.

Artikuluaren luzera beste lau hizkuntzatan.

Honela erabili dugu guk:

- Euskal Wikipediako hutsuneak identifikatzeko.
- Katalanezko Wikipedian Informàtica kategoriako artikuluak, gaztelaniaz eta ingelesez bai baina euskaraz ez zeudenak bilatzeko.



Hobekuntzak Matxin itzultzailean

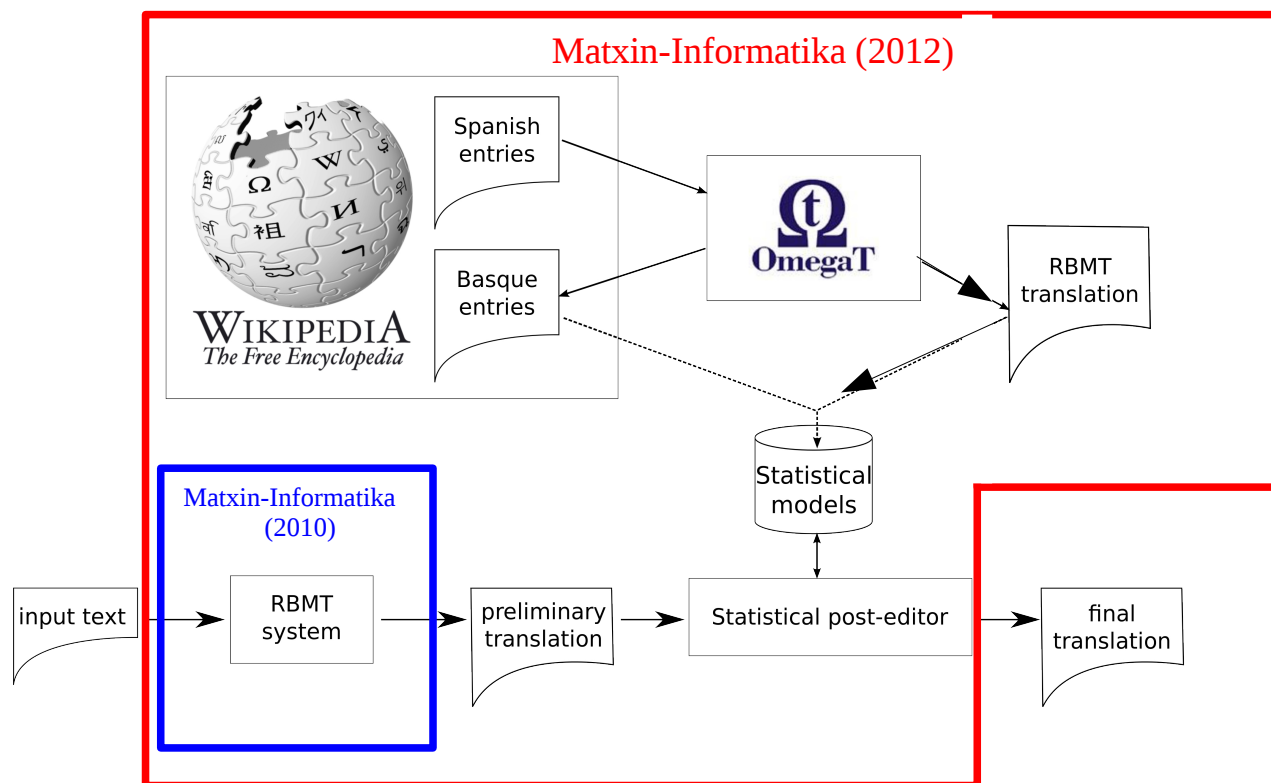
- Matxinen bertsioa informatika gaietarako (2010)
 - Informatikako gaietarako egokitua
 - Hodeian erabiltzekoa (SOAP zerbitzua)
- Espainiera/euskara corpus paralelo bat
 - Mozilla softwarearena. (Elhuyar eta Julen Ruiz).
- Testu itzuliak eta horien eskuzko zuzenketak biltzen dituen corpus bat
 - Espainierazko Wikipediako 100 artikuluko horiek Matxin itzultzailearekin sortutako itzulpenak dituenak, gure kolaboratzaileek egin dituzten zuzenketekin, noski. (50.000 hitz)

Hobekuntzak Matxin itzultzailean (II)

- Matxinen bertsio berri bat postedizio estatistikoaren bitartez hobetua (%10)

'Reciprocal Enrichment between Basque Wikipedia and Machine Translators.'

The People's Web Meets NLP: Collaboratively Constructed Language Resources, Springer, 2013





Etorkizunerako lanak

- Wikipediako estekak eta Wikipediako barne-datuak hobeto erabiltzea postedizioan
 - Itzulpen-sistemaren lexikoa aberasteko
 - Domeinuaren arabera ordain egokiagoak hautatzeko
- Wikipediako informazio hori sakonki erabiltzea
 - Itzulpengintza automatikoan (beste modutara ere)
Proiektu berriak (Tacardi eta QTLeap)
 - Hizkuntzaren prozesamenduan, orokorrean



Etorkizunerako lanak

- Informatika ez den beste arlo batean errepikatzea?

... baina boluntario kopuru minimo bat behar da!

- Boluntario-lana lortzea zaila izan da gurean
 - 100.000 hitzeko itzulpen/zuzenketa egin nahi genuen
 - 50.000 hitzekin moldatu behar izan gara.
- baina azkenean lortu ditugu gure helburuak :-)

OpenMT2 eta Euskal Wikipediatik... **boluntarioei gure eskerrik beroena!**

- Acprisip
- Ana Zelaia Jauregi
- Glabaka
- Irasgo
- Izaskun Etxeberria
- Jiparlaa
- Juanan
- Kgojenola
- Natxo
- Olatzarregi
- Txelo
- UnaiBeleko
- Unai Fdz. de Betoño
- Xsarasola
- Xartola
- Asoraluze
- Anderintxa
- Txibi72
- Joseba Izagirre
- Iñaki.Alegria
- NoraAranberri
- oarbelaitz
- Roberto
- Manex Agirrezabal
- e-gor
- bertolarrieta
- joseba.makazaga
- utolotu
- alexgabi
- Koru
- aingeru
- JesusIbañez
- Jesus Aparizio
- a.illaraza
- ccplaore

