

# Matxin-Informatika:

## versión del traductor Matxin adaptada al dominio de la informática

Iñaki Alegria, Unai Cabezón, Gorka Labaka, Aingeru Mayor, Kepa Sarasola  
IXA Taldea. Euskal Herriko Unibertsitatea  
aingeru.mayor@ehu.es

### Objetivo

Mejorar Matxin, traductor automático es-eu basado en reglas, para el dominio de la informática

### Contexto

Proyecto OpenMT-2  
<http://ixa.si.ehu.es/openmt2>

## Adaptación del léxico al dominio

### • a partir de corpus

#### Recopilación de corpus paralelo

- Del dominio de la informática
- Creado en la localización de Sw

Segmentos	138.000
Palabras es	600M
Palabras eu	440M

#### Tratamiento del corpus

- Analizado, lematizado y procesado con Giza++
- Para cada lema (es) se extraen:
  - sus posibles traducciones (eu)
  - y su probabilidad

#### Uso

- Reordenación de equivalencias en **444** entradas del lexicon  
*dirección: norabide(rumbo) → helbide(ubicación)*

### • a partir de recursos diccionarios

#### Búsqueda en diccionarios

- Accesibles en la red
- Entradas marcadas del dominio de la Informática o similares

#### Uso

- Inclusión de **1623** entradas nuevas en el lexicon  
(sobre todo multi-palabra)  
*base de datos, lenguaje de programación...*
- Modificación de la primera equivalencia léxica en **184** entradas del lexicon  
*rutina: ehitura(hábito) → errutina(procedimiento)*



## Desarrollo de un corpus de postedición para su uso en postedición estadística

### Objetivo

- Creación de un corpus de traducciones es-eu manualmente posteditadas de al menos 100.000 palabras.

### Metodología

- Traducción posteditada de 50 artículos largos de Wikipedia
- En colaboración con la comunidad eu.wikipedia
- Se usa Matxin-Informatika



### Interfaz

- Uso de una Interfaz basada en OmegaT para facilitar el trabajo de postedición
  - Se importa un artículo de es.Wikipedia
  - Traducción automática con Matxin-Inf
  - Corrección manual posteditando
  - Se sube la traducción a eu.Wikipedia



The screenshot shows the OmegaT-2.1.8\_2 interface. The main window is titled "OmegaT-2.1.8\_2 :: probaX". It has a menu bar with "Proiektua", "Editatu", "Ioan", "Ikusi", "Tresnak", "Aukerak", and "Laguntza". The interface is divided into several panes:

- Editorea - Host.UTF8:** Contains the source text: "{referencias}" and a paragraph about "hosts" in computing. A red arrow points to the word "host" in the text.
- Machine Translation:** Shows the translated text: "Host terminoa [[Informatika] informatikan]] konektatuta [[Ordenagailu] ordenagailuei]] [[Konputagailu-sare] sare]] bati kontaktu bere burua erabiltzen dute, ematen duten eta haren zerbitzuak erabiltzen dituzte." Below this is the signature "<Matxin>".
- Parek...:** An empty pane for parallel text.
- Glosarioa:** A pane for the glossary.

At the bottom right, there are status indicators: "0/16 (0/16, 16)" and "179/221".