

# Zapata-kaxatik XML fitxategira euskal hiztegigintzan

ANDONI SAGARNA

---



Hiztegigintzaren hiru urrats nagusiak

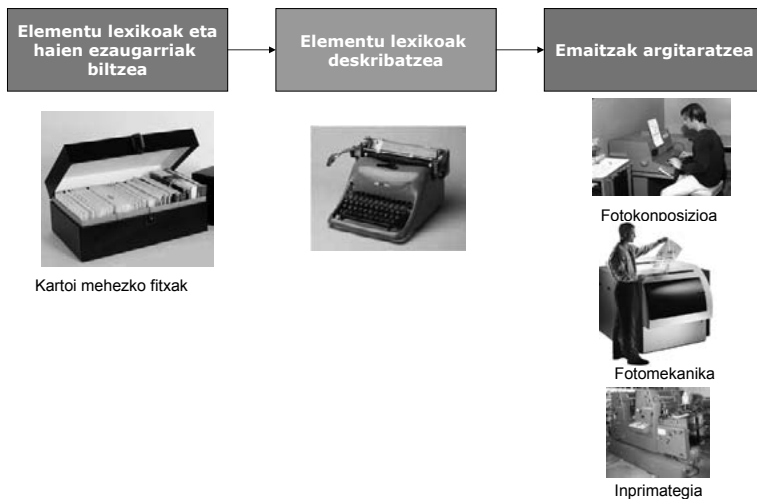
Elementu lexikoak eta  
haien ezaugarriak  
biltzea

Elementu lexikoak  
deskribatzea

Emaitzak argitaratzea

```
graph LR; A[Elementu lexikoak eta haien ezaugarriak biltzea] --> B[Elementu lexikoak deskribatzea]; B --> C[Emaitzak argitaratzea]
```

## Hiztegigintzaren hiru urrats nagusiak 1970eko hamarkadaren lehen urteetan



## Sistema informatiko handiak zituzten enpresa eta erakundeen laguntza 1975-1980

**Unión Farmacéutica  
Guipuzcoana**

1975

Elhuyar Taldeak egindako euskal testuen  
hustuketa euskarri digitalean jartzea.

**EUTG**

1980

**SARASOLA ERRAZKIN, Ibon. 1982.**  
*Gaurko euskara idatziaren maiztasun-hiztegia.*

1977ko euskarazko corpus bat prozesatzea.

**Instituto Deusto**

1980

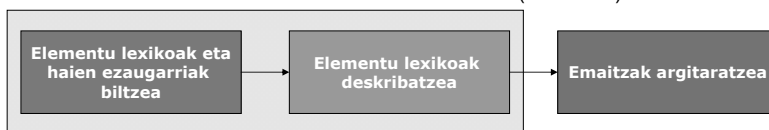
UZEIk egindako terminologia eleaniztunen  
hustuketak euskarri digitalean jartzea.

## UZEIko lehen sistema informatikoaren definizioa 1980

- Enpresetan erabiltzen ziren datu-baseak zurrunka eta garestiak ziren
- Hardware "handia" eta garestia eskatzen zuten
- Informatikariek ezinezkoa ikusten zuten UZEIren lana informatizatzea

### Klaudio Harluxet-en irudimenak ekarri zuen soluzioa

- Propio garatutako datu-base ez-konbentzionala
- Erregistro-egitura librekoa
- Fortran lengoaiari garatua
- Makina "txiki" batek erabil zezakeena (Bull Mini6)



Testu-editoreak oso baldarrak eta garestiak ziren artean

## Orotariko Euskal Hiztegia lantzeko sistema informatikoa 1984 - 2005

*"El autor se convenció de que hoy no se puede hacer una obra de esta naturaleza a base de fichas y cuadernos, como fue el caso de Azkue. Se hacía preciso el recurso al ordenador."*

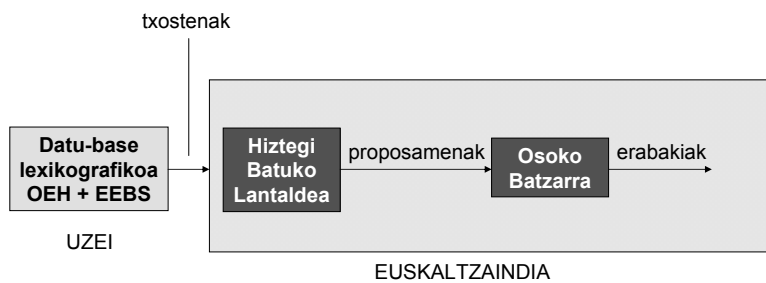
### Aita Villasantek OEHren 1. Liburukiaren aurkezpenean

- Corpusa eta paperezko hiztegiaren originala ordenagailuz prestatzen dira.
- Euskal literaturako 300 obra teklatu bidez digitalizatzen dira
- Hasieran Wang sistema erabiltzen da, gero PCak (1999-2000tik aurrera)
- Corpusa ez da lematizatzen
- Hiztegiaren bertsio elektronikorik ez da argitaratzen

## Egungo Euskararen Bilketa-lan Sistematikoa (EEBS) XX. mendeko corpus estatistikoa 1987 - 2000

- Corpusaren %90 teklatua erabili gabe jartzen da euskarri digitalean
- Datu-base batean antolatzen da
- Corpus osoaren lematizazio erdi-automatikoa egiten da
- Kontsultarako Interneten jartzen da

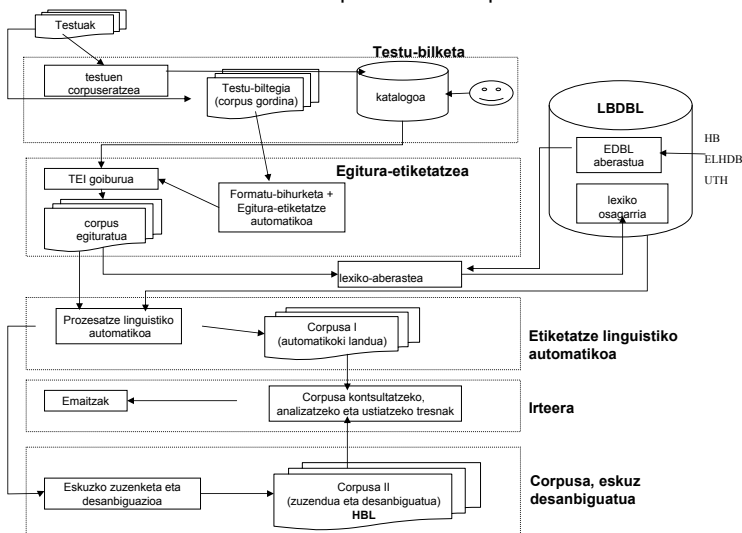
### Hiztegi Batua



**landa-lan**

20.21: **landa-lan** 1, Or *Mi* ("Lan txukunetan etzen ari izaten; gelsez zamarientzat bastak egiten eta landa-lanetan"); **kanpo(-)lan** (eta kanpolar) IE 4: *Mst* ("Gizon pherestu, eta debot batek kanpolar egin behar dutianak, azizenetik bere behizan ikhusten eta aikharen unduan egiteko berheztan dütü"), Etz *Po* ("Lanetan sorthurikan nahiz diren asko, / Ez dire hek eginak lanean hartzeko; / Bulharrak duzte fikro; larrua xurtxo; / Ez dute kuraierik kanpo-lanetako"), Arb *Igand* ("Dakizun bezala, nere auzoko etxean gizon bat bazen Igende-gabea, Meza-Bezperetako denboran kanpo lanetan trebeki ari zena"), JE *Bur* ("egun osoa barne ala kanpo lanetan iragan ondoan, jendeak, kabalak, gauzak zoin beren lekuean, tenoretan eta izarian moldaturik, azkenik zotan behi oherat"), eta OEH argitaratua, Prop 1906 ("Haurek egiteko dutena da, baratze eta alhor, hori da kanpo-lanetan, gure laguntza") eta alp. *Aran-Bago MarkWed*.  
 20.22: a) *Ikerketa-motari dagozkionak dira: landa-lan 1, EHHistoriaz* ("Data hau euskal azarrietatik J. M. Barandiaranen landa-lanetan jardun ziren ikertzaile gazteen Doktoratze-tesi edo lan espezializatuaren argitalpenetan oinarritzen da"), eta **kanpo-lan** 6: I. Inurrieta 5 (adib., "kanpo-lana, (...) ikasleak hirira irango dira", "Utzi ikasleei beren laguntzaileak gelara ekartzan, kanpo-lana berregitiko", "Prestatuta. Kanpo-lanari ekin baino egun batzu lehenago, ikasleek entzun edo irakurritako bixteak eramango dituzte ikasgelara"), Bertsolaritze ("Ikasuntu kultural batetik gaurko bertsolaritzaren berballoketak, berari buruzko kanpo-lan zabal bat egin ahal dadin giroa eta materiala prestatzen ditu"); b) cf. E. J. Zubillaga ("Emeretzigarren eunikian [siglo 19] bere erdialdetik gora Argentina'ra etorritako euskaldunak orderaz mintzitzen ez zekirenak, la kanpo lanetan jarduntzen ziranak beren alaba-semeak euskeraz ikasi zuten apur bat, beren aita-amak beti euskeraz mintzaten ziralako"), UZEI ("Asko erabiltzen da kanpo-lanetarako eta gure artean aitzainuz ere bai").  
 20.20: **landa-lan**. Euskalterm 2; **kanpo-lan**. Euskalterm 2 // Ez dugu aurkitu ap. AB38, AB50, DFrec, HiztEn, LuE // *Eta trabajo de campo* itzuliak: **eremu-ikerketa**: Euskalterm 1; **ikerzango**: AB38 1; **lan bortzatuaren esparru**: Euskalterm 1; **lan sail**: AB38 1; **landa-azterketa**: Euskalterm 1; **laneko udaleku**: Euskalterm 1 (campos de trabajo); **tokiatokiko lan**: AB50 1.  
 20.81: **kanpolar**: HiruMila (trabajo de fuera), Elhiz (trabajo de fuera, trabajo de campo), Casve *EF* (travail extérieur), HaizeG *BF* (travaux des champs), LH *DBF* (travaux du dehors et plus communément travaux des champs), DRA (**kanpo-lan**: trabajo de fuera, y más generalmente trabajo del campo) // Ez dugu aurkitu ap. EuskHizt, EskolaHE, Lur *EG/CE* eta *EF/FE*, XarHiz, PMuj *DVC*.  
 20.83: Erdial *trabajo de campo* // *travail sur le terrain*; *travaux des champs* formen ordainak: **landa-lan**: HiruMila, Elhiz // HaizeG *FB*: Ø / *lur-lan*, **landa-lan** // Ez dugu aurkitu ap. Lur *EG/CE* eta *EF/FE*, XarHiz, Casve *FE*, T-L *LBF*, PMuj *DCV*, 20 8 4 (DLF); Ø (k. Le Petit Robert, adibideetan. *Faire une enquête sur place*); IT (S. Carbonell); Ø; en (Collins); *fieldwork*; **de** (Langenscheidts); Ø.

**Lexikoaren Behatokia: Corpora lantzeko prozesua**





#### **Egitura-etiketetzea.**

TEI P5 (XML).

**Automatikoa:** testuaren egitura-ezaugarriak (atalburuak, atalak, azpiatalak, paragrafoak, zerrendak, taulak, oin-oharrak, irudi-oinak, eta abar.); nabarmentze-ezaugarriak (tipografikoak, hau da, letra-estiloa eta komatxo edo kakotxak)

**Eskuz lantzekoa:** nabarmentze-ezaugarrien balioa (enfasia, aipua, atzerri-hitza, ohiz kanpoko adiera edo erabilera ironikoa, metahizkuntza, terminoa, izen berezia...)



#### **Aurreprozesamendu linguistikoa**

Corpusaren lexikoi osagarria elikatzea (EDBLn ez dauden eta corpusean atzeman diren maiztasun handiko lemak)

Aldaeren normalizazioa (<reg> etiketa)

Akats ortotipografikoen zuzenketa (<corr> etiketa)



## Analisi linguistikoa

Lema, kategoria, azpikategoria, kasua

