

Web-corpusen Ataria



Igor Leturia, Iñaki San Vicente, Iker
Manterola, Antton Gurrutxaga

IEB 2013

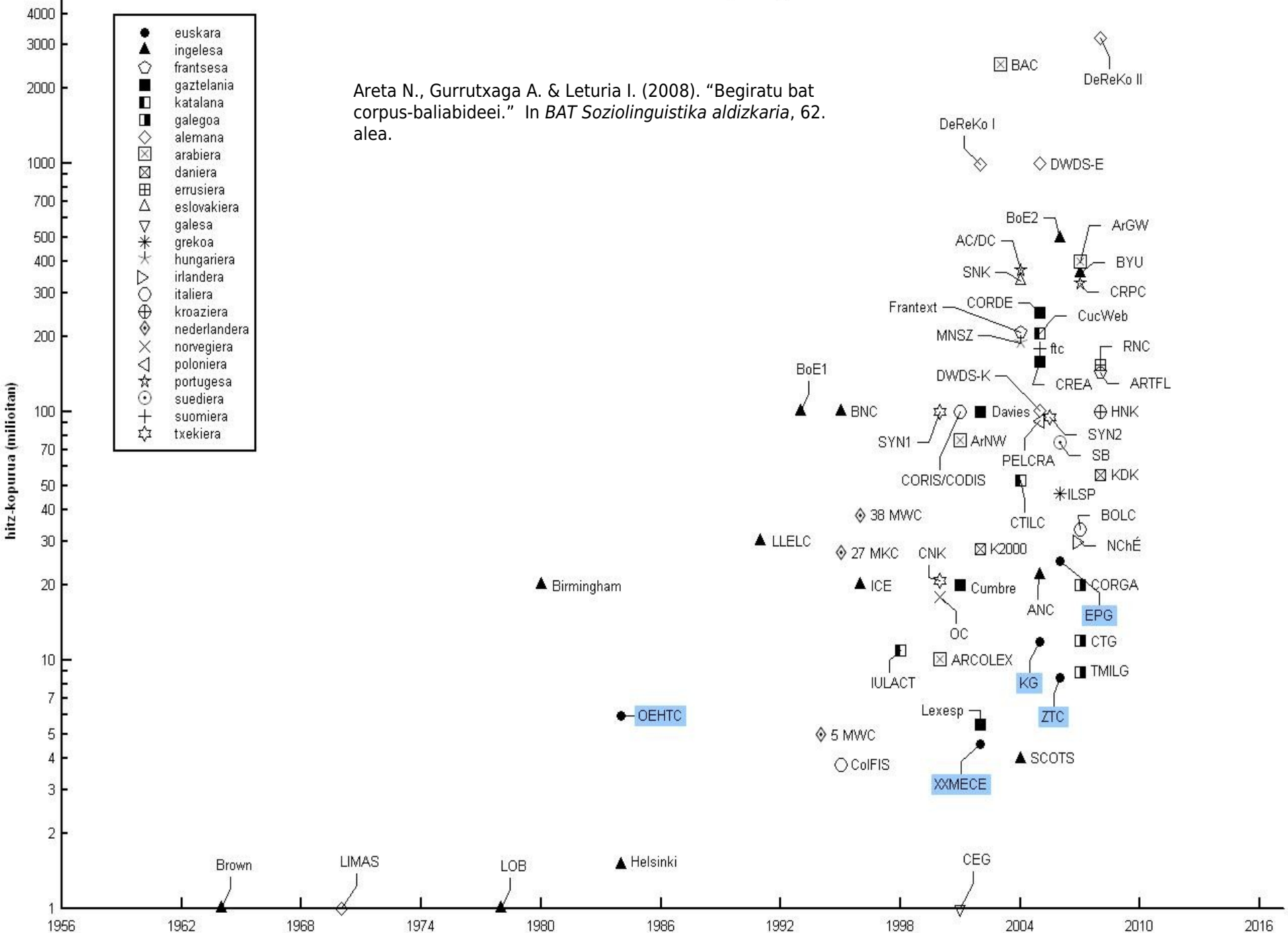
Miramon, 2013ko maiatzaren 8a

Edukia

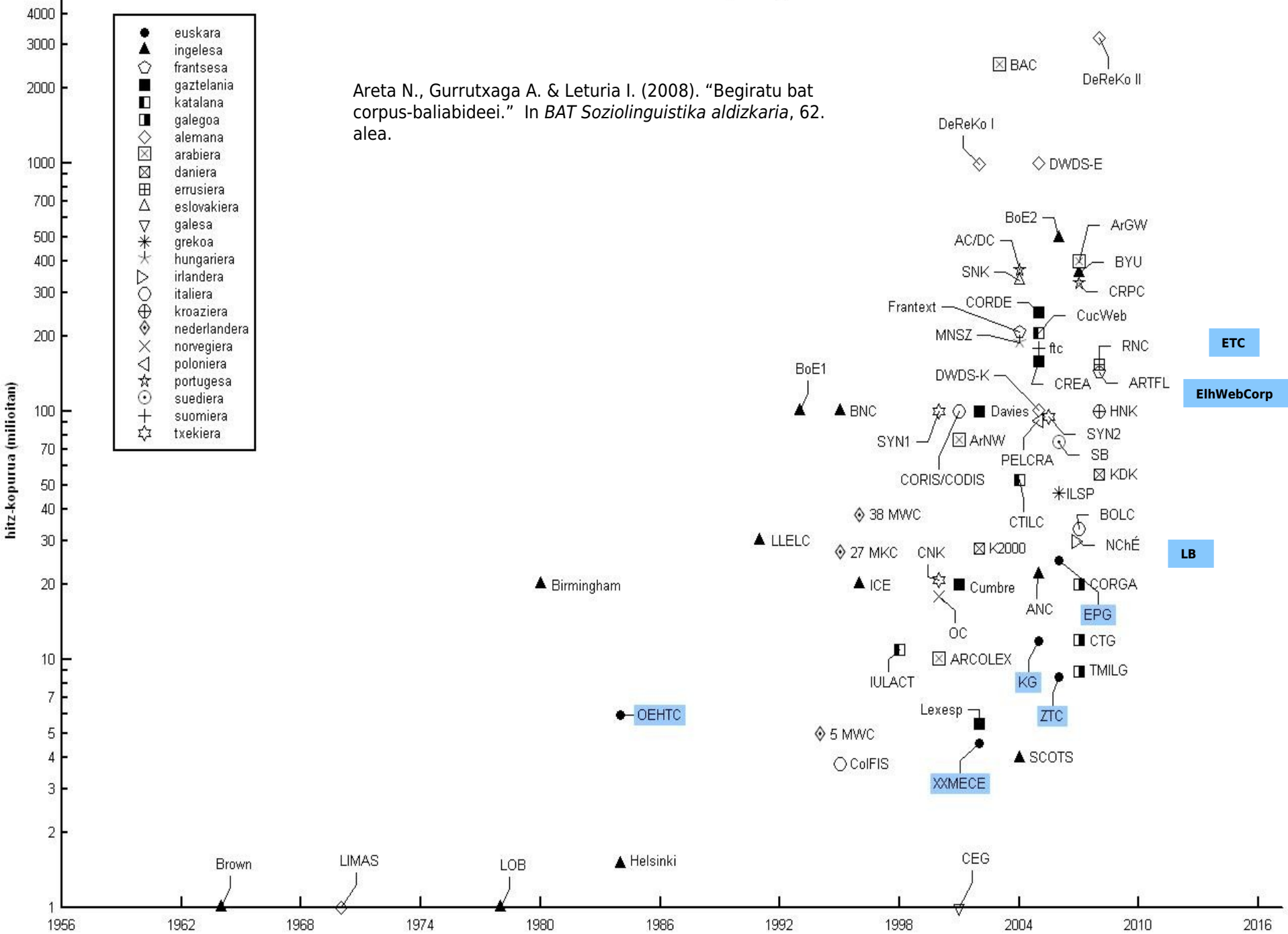
- ZERGATIK egin dugu?
- ZER da?
- ZENBAT du?
- NOLA egin dugu?
- ZERTARAKO balio dezake?
- Hurrengo urratsak

ZERGATIK egin dugu ?

Zenbait hizkuntzatako testu-corpus nagusiak



Zenbait hizkuntzatako testu-corpus nagusiak



Areta N., Gurrutxaga A. & Leturia I. (2008). "Begiratu bat corpus-baliabideei." In *BAT Soziolinguistika aldizkaria*, 62. alea.

ZERGATIK egin dugu?

- Ikuspegi “praktikoaz” gain...
 - Ikuspegi linguistikoa
 - Internet ezin ukatuzko errealitate "linguistikoa" ere bada
 - Interneten BAKARRIK argitaratzen diren testuak gero eta ugariagoak dira, eta ezaugarri bereziak dituzte
 - Interneteko testu asko ez dira "diskurtso kontrolatuak", bat-batekoak baizik; hizkuntzaren erabilera errealia islatzen dute
- Interesgarria da webaren alderdi linguistiko bereziak aztertzea

ZER da?

- Zer da Web-corpusen Ataria?
 - Webetik automatikoki osatutako euskarazko corpusen kontsulta aurreratua
 - Elhuyar Fundazioko hizkuntza-teknologien I+G taldeak egindako lanaren emaitza
 - Hiru atal:
 - Euskarazko web corpus elebakarra
 - Euskara-Gaztelania web corpus paraleloa
 - Hitz-konbinazioen kontsulta

Galdera

Zer	Aukerak	Bilatu	Kategoria
Lema	Da	xedapen	
Dist.	Non	Zer	Aukerak Bilatu
Dist.	Non	Zer	Aukerak Bilatu

Emaiza
 Testuinguruak eta kopuruak

Ordenatu honen arabera
 Dokumentua

Kopuruak **Gehienez** %

Forma	10	
Lema		
Kategoria		
Lema eta kategoria		

➔ Bilatu ✖ Garbitu

Bilaketa arrunta

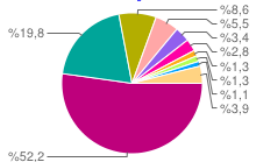
Emaitzak: 13486

Kopuruak

Denak (13486)

Forma	Kop
xedapen	7041
xedapenAK	2668
xedapena	1158
xedapenen	736
xedapenetan	463
xedapenek	380
xedapenei	181
xedapenean	176
xedapenetatik	151
Beste guztiak	532
Guztira	13486

Guztien testuinguruak batera



- https://www6.euskadi.net/u81-0003/eu/contenidos/autorizacion/2216/eu_4901/eu_17719.html (1)
- ...datzen ere prozedura administratiboak erregulatzen dituzten **xedapen** arautzaileak, eta, ondorioz, ezin ditu sortu berariazko Ara...
- http://www.kubxa.net/wkn_corporativo/2009/informe_financiero/eu/memoria_consolidada_02.htm (1)
- ...kal, Administrazio eta Ordena Sozialari buruzkoaren, azken **xedapenetako** hamakaigarrenak ezartzen duenaren arabera.
- http://www.euskadi.net/cgi-bin_k54/ver_e7CMDVERDOC8&ASEB03D&DOCN0000381018CONF/config/k54/bopv_e.cnf (2)
- ...3. artikuluetan eta artikuluekin erlacionatutako beste **xedapenetan** ezarritakoarekin bat etorri.
- AZKEN XEDAPENAK
- <https://ssl4.gipuzkoa.net/euskera/gao/2004/02/16/e0401200.htm> (4)
- XEDAPEN GEHIGARRIAK
XEDAPEN IRAGANKORRA
AZKEN XEDAPENAK
I. KAPITULUA. BABES-ARAUBIDEAREN XEDAPEN OROKORRAK
- <http://www.cgt-lkn.org/laescoba/convenio/Mungia.pdf> (3)
- ...a jatorria edozein izanik ere. Aplikazio orokorreko legezko **xedapenak**, edo Hitzarmen Kolektiboak, edo banakako kontratuak direla ...
- XEDAPENAK
40. artikulua.--Azken **xedapenak** Hitzarmen hau indarrean dagoen bitartean, bertan itundutako...
- <http://www.senado.es/brsweb/CALEX/textos/paisvasco/44/2003/03.pdf> (7)
- Bestelako **Xedapenak**
Xedapen Orokorrak
- ... hogoita bat artikulua ditu, lau **xedapen** iraqankor eta azken **xedapen** bat. Legeak zehaztu egiten ditu...
- Aurrezki-kutxek lege honetako **xedapenetara** eaqokitu beharko dituzte beren estatutuak eta araudiak, eta ...
XEDAPEN
XEDAPENA
- <http://www.innobasque.com/Modulos/Documentos/Visor.aspx?docId233> (3)
- Xedapen** Orokorrak
- ...azaroaren 26ko 30/1992 Legean zehaztutakoa edo bere ordezko **xedapenak** ezartzen duena izanqo da ezarzarri. Hirugarrena.Å Aurrekon...
- ...zarritako mugen barruan. Aipatutako Lege Organiko horren 5. **xedapen** gehigarrian ezarritakoaren arabera, Unibertsitatea Koordina...
- <http://www.donostia.org/ordenanzas.nsf/dbda625a50d6b4b8412566800044e4ed/b8970ea57f9ec357c12574ea002a49bc?OpenDocument> (2)
- ...Iaren Seko Foru Arauean ezarritako prezio publikoari buruzko **xedapenekin** bat, Donostiako Gazteriaren Aterpetxe eta Kanpalekuen Udal ...
10. art.- Azken **xedapena**
- <http://www.euskadi.net/bopv2/datos/1993/09/9302983a.pdf> (1)
- Bestelako **Xedapenak**
- <http://www.gipuzkoa.net/euskera/gao/2004/12/03/e0411069.htm> (5)
- XEDAPEN OROKORRAK
Xedapen Gehigarria.
Indargabetzeko **Xedapena**.
Azkenurreko **Xedapena**.
Azken **Xedapena**.
- http://www.euskara.euskadi.net/r59-uhadm2/eu/contenidos/informacion/ohiko_galderak/eu_irati/dekretuari_buruzkoak.html (2)
- ...tik kentzeko espedientea has daitekeela ondoriozta daiteke. **Xedapen** honen oinarria 6/1989 Legean daqo, 50. artikuluan hain zuze...

Galdera

Zer **Aukerak** Bilatu
Lema Da

Dist. Non Zer Aul

Dist. Non Zer Aul

Bilaketa arrunta

Emitzak: 13486

Kopuruak

Forma

Forma

- xedapen
- xedapenAK
- xedapena
- xedapenen
- xedapenetan
- xedapenek
- xedapenei
- xedapenean
- xedapenetatik
- Beste guztiak

Guztira

Guztien testuinguruak ba

Web-corpusen Ataria
Corpus paraleloa

Galdera

Hizkuntza Zer bilatu
Euskara Lema

Hizkuntza Zer bilatu
Gaztelania Lema

Emitzak: 81 itzulpen-unit

Ukidezia

Ukidezia

CAN Goreneko **Likidezia**

Behin pizgarria lortu denaz **likidezia** duen produktu bar aurrezten jarraitu nahi baduzu, une bakoitzean ego

Berankortasun maila murriztu

Ondasuna eskura dago norb ez du **likideziarik** galtzen ez batzuetara bideratzean.

Zerga abantaila maximoekin gero.

Aurreikuspen plana (EPSV) i betetzen denean.Pentsio p ezintasunagatik ordaintzen e hiltzen bada, beste pertsona:

Banca Cívica taldea izan finantza-sistema jasaten an sektorean -lehendabizikoa eta egoera sendoa du, **likideziagatik**

2011ko ekainaren 21ean, **likidezia** berretsi dio kutxi **likideziagatik**

Likidezia: zure dirua (osoril duzu, inolako gainordainik e

Erabateko **likidezia** eskaintzi aukera baitago.

Inbertsio fondo bat kontrat izateaz gain, **likidezia** da, erabilgarri izateko aukera en

Berehalako **likidezia**, izatek utzita. Epe laburrean amaitz

Aseguru hori gengarria da izateko: **likidezia**, gaudimen

CAN GORENEKO **Likidezia** epe labur, ertain eta luzerak CANek CAN GORENEKO **Likidezia** aktiboekin **likidezia** sortzeko

Web-corpusen Ataria
Hitz-konbinazioak

Galdera

1. lema **Konbinazioak** **Ordenatu honen arabera**

2. lema

t neurria LLR PMI PMI³ x² Fisher

Galdera

Konbinazioa	f	f1	f2	t neurria	Adibideak
energia aurreztu	825	3158	2120	28,68	Iraunkorra ez den egoera hori bideratzeko asmoz, energia aurrezteko planak eta energia berriztagarriak sustatzeko laguntzak baliatu nahi izan ditu Europako batasunak; lobby nuklearak, aldiz, are iragangaitzago egin nahi du egoera, erreaktore nuklear berriak saltzeko asmoz. adib. --
energia sortu	386	3158	64146	17,9	Halaber, hidrokarburuak neurrigabe erretzen dira energia sortzeko . adib. --
energia lortu	285	3158	44605	15,47	Babidik borroka honetan jasandako zaurietatik Boo ahaltsua esnatzeko energia lortu zuen. adib. --
energia kontsumitu	167	3158	1348	12,87	energia kontsumituko den tokian ezartuta izan daitekeen energi-iturri moderno bat da. adib. --
energia ekoiztu	116	3158	1373	10,7	Gure gizarteak energia ekoizteko sekulako baliabide kantitatea behar du, batez ere erreagai fosiletan oinarritzen delako. adib. --
energia erabili	176	3158	72806	10,32	Etkeak, patioan kokatuta, eguzkia, haizea eta abarren energia erabiliz zeinbat gauza egin zezakeen erakusten zigun. adib. --
energia aprobetxatu	76	3158	1110	8,65	Hezitzailearen egitekoa da, jolasean askatutako energia aprobetxatzea , balioetan eta bizi egitasmoetan heztzea eraginez. adib. --
energia askatu	72	3158	1561	8,39	Lur barnean zeuden substantzia erradioaktiboen desintegrazioan askatutako energia sumendi-erupzioetan askatutako energia . adib. --
energiatz bete	73	171	49414	8,38	Eta hori modu enpatikoan egiten zuten: hitzak gutxieneko adierazpena zeuden murriztuta, bizitza modern o a ren izaera paradoxikoa erakutsiz, hau da, teknologien garapenaren ospakizuna eta, era berean, alienazio-sentimendu sendoa teknologia haien azpian izateagatik: gure batera aldatzen ari gara eta orain energiatz betek gaude. adib. --
energia metatu	63	3158	423	7,91	Kondentsadore batek U energia metatzen du, eremu elektriko gisara, eta geroago frogatuko dugunez: kondentsadore lau eta paraleloa hasteko, kalkulatu dezagun plano mugagabe batek sortutako eremu elektrikoak, s karga-dentsitatea duela kontsideratuz, eta Gauss-en legea erabiliz. adib. --
energiatz hornitu	60	171	1264	7,74	Makineria guztia energiatz hornitzeko , 5, 5 megawatteko zentral elektriko bat eraiki da lurpean; hiru sorgailuz osatutakoa. adib. --
energia galdu	77	3158	18024	7,67	Guk, esate baterako, zera ikertzen dugu, alegia, elektroiak zenbat denbora irauten duen egoera aldatzen zaionean, atzera energia galdu eta beste egoera batera igaro baino lehen. adib. --
energia berreskuratuz	51	3158	3489	6,88	Hiri-hondakin solidoetatik energia berreskuratuzeko plantak. adib. --
energia bihurtu	59	3158	12926	6,78	Eguzkiarekiko gure estimuak haren erregetza luzatu nahi du beharbada, eta argi bihurtutako energia itzuli nahi diogu eguzkiari. adib. --
energia gastatu	43	3158	970	6,48	Eta nire harridura handitzen joan da; izan ere, non egon gara euskaldunok, non egon dira euskal ikerlariak orain arte? kanpotik etorri behar izan behar du katalan batek honen guztiaren berri emateko? zenbat energia gastatzen dugun exigentzia, eskari eta protestatan, berauek oinarri eta indar litzakeen ikerketetarako erabili beharrean imagina genezake nafar ikerlari bat el Vendrell-eko katalanaren garapena aztertzen? Ruper Ordorikak kontzertu zoragarria eskaini digu duela gutxi Lezaman; lagunartekoa. adib. --
energiatz baliatu	38	171	4237	6,14	Argiari Lasarte-Orian duen pabiloian eguzki energiatz baliatzen zen sistema bat instalatzeko zer eskatuko liokeen azter-tzea otu zitzaion. adib. --
energia xahutu	33	3158	339	5,71	Marruskadurak energia xahutzen du, energia zinetikoa bero bihurtuz. adib. --
energia ustiatu	30	3158	495	5,43	Biomasa tratatzeko makineria espezifikoa, energia ustiatzeko landan, bilketa eta garraioa errazteko, eta hartara, biomazaren garraioari lotutako gastuak gutxituko dira. adib. --
energiari egon	27	27	229583	4,99	adib. --
energiari egon	27	27	229583	4,99	Izan ere, muskulatura, mugitzeko motorra izateaz gain, karozeria ere bada; eta, kasu honetan, ez da komeni karozeria astun eta garestirik energiari dankonez . adib. --

Web-corpusen Ataria
Corpus paraleloa

Galdera

Hizkuntza Zer bilatu
Euskara Lema

Hizkuntza Zer bilatu
Gaztelania Lema

Emitzak: 81 itzulpen-unit

Galdera

1. lema **Konbinazioak** **Ordenatu honen arabera**

2. lema

t neurria LLR PMI PMI³ x² Fisher

Web-corpusen Ataria
Hitz-konbinazioak

Galdera

Hizkuntza Zer bilatu
Euskara Lema

Hizkuntza Zer bilatu
Gaztelania Lema

Emitzak: 81 itzulpen-unit

Galdera

1. lema **Konbinazioak** **Ordenatu honen arabera**

2. lema

t neurria LLR PMI PMI³ x² Fisher

Web-corpusen Ataria
Corpus paraleloa

Galdera

Hizkuntza Zer bilatu
Euskara Lema

Hizkuntza Zer bilatu
Gaztelania Lema

Emitzak: 81 itzulpen-unit

Galdera

1. lema **Konbinazioak** **Ordenatu honen arabera**

2. lema

t neurria LLR PMI PMI³ x² Fisher

ZENBAT du?

- Zenbat du euskarazko web-corpus elebakarrak?
 - Hitzak: 125 milioi inguru
 - Dokumentuak: 82 mila inguru
 - Domeinu ezberdinak: 6.206
- Testu-motak: periodistikoa, literatura, akademikoa, hezkuntza, webeko berezkoak (blogak, foroak...), administratiboa...

Domeinua	Dokumentuak	Hitzak
www.argia.com	2.886	4.026.006
www.zientzia.net	2.353	3.093.083
www.euskomedia.org	757	2.638.731
eibar.org	1.299	2.538.105
eu.wikipedia.org	2.422	2.281.205
www.soziolinguistika.org	690	2.102.094
www.ikasbil.net	1.220	1.932.899
klasikoak.armiarma.com	905	1.930.030
www.uztaro.com	407	1.916.843
www.euskonews.com	1.354	1.743.869
www.euskadi.net	775	1.683.893
www.armiarma.com	600	1.601.705
www.jakingunea.com	581	1.583.140
revista.consumer.es	1.518	1.459.116
paperekoa.berria.info	1.769	1.448.311
www.euskaltzaindia.net	563	1.388.255

ZENBAT du?

- Zenbat du euskara-gaztelania web-corpus paraleloak?
 - Hitzak: 18 milioi inguru
 - Segmentuak: 1,2 milioi inguru
 - Dokumentuak: 174 mila inguru
 - Domeinuak: 84
- Testu-motak: administratiboak, erakundeen webguneak, turismoa, kultura, albisteak, enpresak, alderdi politikoak...

Domeinua	Dokumentuak	Segmentuak	Hitzak
www.euskonews.com	3.444	165.736	3.014.432
www.euskalkultura.com	9.888	98.042	1.679.440
www.hiru.com	20.116	104.738	1.418.600
www.donostia.org	8.594	84.502	1.379.709
www.euskomedia.org	1.356	41.998	998.369
www.ermua.es	12.930	66.886	971.038
www.3digitala.com	16.888	66.028	829.583
www.turismoa.euskadi.net	4.528	41.530	513.617
web.bizkaia.net	3.946	31.126	398.359
www.euskadi.net	2.884	28.604	355.664
www.eusko-ikaskuntza.org	4.426	16.544	306.171
www.ibaizabal.com	5.252	27.820	305.847
www.pasaia.net	1.716	22.664	291.672
www.mondragon.edu	2.772	20.812	286.619
www.lezo.net	892	15.598	279.035
www.euskara.euskadi.net	898	16.584	256.107
www.catedralvitoria.com	2.548	14.040	251.208
www.aralar.net	1.458	12.124	247.887
www.guggenheim-bilbao.es	3.036	17.304	235.906
www.basqueresearch.com	3.836	15.980	234.674

ZENBAT du?

- Zenbat du konbinazioen atalak?
 - Konbinazioak: 133 mila inguru
 - Agerpenak: 10 milioi inguru

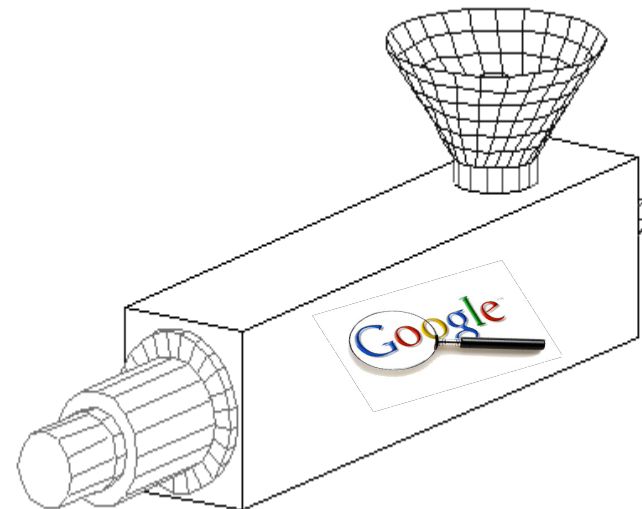
Mota	Kopurua	1. osagaiak	2. osagaiak
IZE-ADI	71.550	8.199	1.528
IZE-ADJ	27.931	4.797	2.782
IZE-IZE	33.759	7.067	7.002

NOLA egin dugu euskarazkoa?

XX. Mendeko Corpus Estatistikoko maiztasun handieneko 2.000 hitzak

Ausazko bikote-konbinazioak bidali web-bilatzaile baten APIari, sorkuntza morfologiko bidezko hedapena eginez, eta hizkuntza iragazteko hitzekin batera

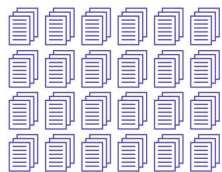
Euskarazko web-corpus gordin erredundantea



NOLA egin dugu euskarazkoa?



Euskarazko web-corpus garbia



Eustagger

Euskarazko web-corpus etiketatua

Elhuyarren CAS → Web-corpusen Ataria

- **Luzera-iragazkia:** laburregiak eta luzeegiak iragazi, errore-orriak, testu gutxi dutenak eta spama saihesteko
- **Hizkuntza-iragazkia:** Langld hizkuntza-detektatzailea paragrafo mailan, testu elebidunetako zati ez euskarazkoak kentzeko
- **Orrien garbiketa:** nabigazio-menuak, oin-oharrak... kentzeko
- **Errepikatuen iragazkia:** ia berdinak diren dokumentuak iragazteko (agentzietako berriak, lizentzia libredun testuen kopiak...)
- **Barne-hartze iragazkia:** blog edo komunikabide bateko artikulua bat osorik orri nagusian agertzen denean bi bider ez sartzeko

NOLA egin dugu paraleloa?

Euskarazko hitz-konbinazioak



Euskarazko edukia duten domeinuak:

www.elhuyar.org

www.zientzia.net

Eduki elebiduna duten domeinuak



eu



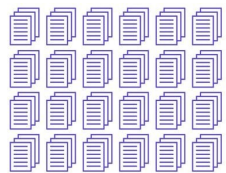
es



Domeinuak eduki elebiduna du?
Alt testuak,
Hizkuntza arteko estekak,
Combo egiturak,...

NOLA egin dugu hitz-konbinazioen erauzketa?

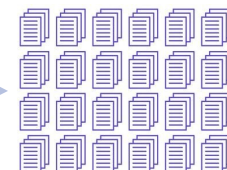
Euskarazko web-corpora etiketatua



- Corpus elebkar etiketatuaren irteera moldatu
- Interesatzen zaigun informazioa atxiki
- Anlisi batzuk testuinguruaren arabera eraldatu

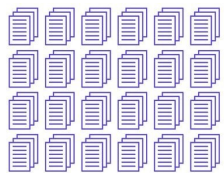
- egunean egun IZE ARR INE NUMS
- dokumentatutako dokumentatu **ADI_SINPART_GEL_0** grafien_grafia_IZE_ARR_GEN_NUMP
- indar_indar_IZE_0_0 armatuak_armatu_ADI_SINPART_ABS_NUMP → indar_indar_IZE_0_0 armatuak_armatu_ **ADJ_ADIPART_ABS_NUMP**
- ...

Bigrama-erauzketarako web-corpora etiketatua



NOLA egin dugu hitz-konbinazioen erauzketa?

Bigrama-erauzketarako web-corpus etiketatua



Patroi morfosintaktiko batzuk betetzen dituzten IZEADI, IZEADJ, IZEIZE agerkidetzetatik bigramak sortu ($w = \pm 1$; $f > 10$)

- IZEADI: izenaren forma_lema_kat_azpik_kasua_mugat; aditzaren lema (+aspektua, flexioa...)
- Izenaren formak normalizatu, f handienekora
- IZEIADJ: iz.Ø+izond. sintagmak edo sintagma-zatiak
- IZEIZE: iz. Ø +iz. sintagmak edo sintagma-zatiak (elkartek)

Bigrama-erauzketa



Bigrama-erauzketa + informazio estatistikoa

id	l1	l2	f	f1	f2	N	am.t.score	am.log.likelihood	am.MI	am.MI3	am.chi.squared	am.z.score
10443	nahi_nahi_IZE_ARR_ABS_MG	ukan_ADI	115928	143684	426774	2817230	276.5541924	333326.1246	0.72640428	10.85478097	505872.384	638.237604
2328	behar_behar_IZE_ARR_ABS_MG	ukan_ADI	102522	214368	426774	2817230	218.7697593	142397.6740	0.49928146	10.52091560	192738.624	388.712066
13257	uste_uste_IZE_ARR_ABS_MG	ukan_ADI	56044	56371	426774	2817230	200.664221	214195.3541	0.81709774	10.31415599	317803.094	514.066433
2315	behar_behar_IZE_ARR_ABS_MG	izan_ADI	99512	214368	629138	2817230	163.6992368	66382.3628	0.31779189	10.31354280	77627.050	236.016626
11169	parte_parte_IZE_ARR_ABS_MG	hartu_ADI	18292	22636	93411	2817230	129.6985361	106445.8006	1.38688579	9.91140818	427467.972	640.291844
8823	kontuan_kontu_IZE_ARR_INE_NUMS	hartu_ADI	14494	19298	93411	2817230	115.0761443	79664.0841	1.35510024	9.67747675	312392.186	547.690769
226	ahal_ahal_IZE_ARR_ABS_MG	izan_ADI	22077	27128	629138	2817230	107.8104326	43202.0911	0.56159867	9.24947878	55865.450	205.806991
5239	ezin_ezin_IZE_ARR_ABS_MG	izan_ADI	26592	48587	629138	2817230	96.5327738	24447.4904	0.38930730	9.23880931	29920.507	151.122259
9137	lan_lan_IZE_ARR_ABS_MG	egin_ADI	18584	55657	300597	2817230	92.7606719	21256.9122	0.49545707	9.03373546	30750.560	164.093893
2365	behera_behe_IZE_ARR_ALA_NUMS	utzi_ADI	8278	16263	28987	2817230	89.1443575	55953.1793	1.69434379	9.53019464	399519.065	626.997244
3178	deta_deta_IZE_ARR_ABS_NUMS	egin_ADI	9791	11233	300597	2817230	86.8366814	35812.4581	0.91216895	8.89382306	69232.961	248.191952
6904	hitz_hitz_IZE_ARR_ABS_MG	egin_ADI	11337	20466	300597	2817230	85.9662667	24982.5302	0.71530284	8.82429913	43264.120	195.875215
4103	erabakia_erabaki_IZE_ARR_ABS_NUMS	hartu_ADI	7961	10511	93411	2817230	85.3184828	43424.5220	1.35874787	9.16068312	172623.824	407.771297
5248	ezin_ezin_IZE_ARR_ABS_MG	ukan_ADI	19090	48587	426774	2817230	84.8953307	17026.5353	0.41391020	8.97552205	22416.856	136.721257
9036	lagun_lagun_IZE_ARR_ABS_MG	hil_ADI	6423	17523	21605	2817230	78.4668589	41778.0415	1.67939749	9.29487334	298416.712	542.481604
4822	espero_espero_IZE_ARR_ABS_MG	ukan_ADI	9651	13323	426774	2817230	77.6951952	22108.4695	0.67959464	8.64873927	34180.878	169.899258
5308	falta_falta_IZE_ARR_ABS_MG	izan_ADI	10739	12829	629138	2817230	75.9830069	21979.0205	0.57384742	8.63577511	27991.057	147.109320
10516	neurriak_neurri_IZE_ARR_ABS_NUMP	hartu_ADI	6364	9912	93411	2817230	75.6549226	31080.8866	1.28699316	8.89445350	115038.258	332.915313
1664	aurrena_aurre_IZE_ARR_ALA_NUMS	egin_ADI	8739	17482	300597	2817230	73.5289580	17027.2496	0.67070829	8.55363176	28531.913	159.152166
9927	martxa_martxa_IZE_ARR_INE_NUMS	Jarri_ADI	5316	7367	47148	2817230	71.2199219	35447.6728	1.63465333	9.08582328	223009.393	467.657775
2422	berri_berri_IZE_ARR_ABS_MG	emoi_ADI	5480	10005	110195	2817230	68.7405445	22355.8524	1.14622388	8.62378500	69107.205	257.231940
11221	partida_partida_IZE_ARR_ABS_MG	jokatu_ADI	5036	14837	36012	2817230	68.2922156	68.1841466	1.42411469	8.82828613	126106.546	351.907519
341	akordioa_akordio_IZE_ARR_ABS_NUMS	lortu_ADI	4921	11638	33380	2817230	68.1841466	28691.4956	1.55251098	8.93661771	168597.863	407.328270
1730			1730	1730	1730	1730						

ZERTARAKO balio dezake?

- <http://webcorpusak.elhuyar.org>

- Euskarazko web-corpus elebakarrean
 - Handitu crawling metodoa erabiliz (bilatzaileetatik ezin da askoz gehiago lortu) eta online jarri
 - Corpuseko testuak karakterizatu (eremua, generoa...)
- Euskara-gaztelania web-corpus paraleloan
 - Esalditik beherako unitate-mailan parekatu (hitzak, terminoak, fraseologia...)
- Hitz-konbinazioak
 - Idiomatikotasuna karakterizatzeko teknika aurreratuak (konposizionaltasuna)
 - Konbinazio-mota gehiago: IZLG+IZE, ADB+ADI...
 - Patroi morfosintaktikoak (argumentuak, azpikategoria...)
 - Fraseologia elebiduna

Web-corpusen Ataria

Eskerrik asko!

Web-corpusen Ataria



Igor Leturia, Iñaki San Vicente, Iker
Manterola, Antton Gurrutxaga

IEB 2013

Miramon, 2013ko maiatzaren 8a