

# Itzulpen-corpusak hiztegitintzan: baliabide eta bide berriak

DAVID LINDEMANN

Itzultzailea eta Euskal Herriko Unibertsitatean doktoregiaia

## 1. Sarrera

Egun, inork ez du zalantzan jartzen corpus hizkuntzalaritzak nabarmen aberastu duela hiztegitintza (ikus gaiaren oinarritzko aurkezpenak Atkins & Rundell 2008; Svensén 2009). *Corpus revolution* deitutako berrikuntzari esker (Rundell & Stock 1992; Krishnamurthy 2002), hizkuntzaren erabilera hobeto deskriba daiteke gaur. Corpus garaiaren aurreko hiztegiek, zenbait egilek nabarmendu dutenez (Klosa 2007), hitzen adiera (eta hiztegi elebidunez ari garenean, adieraren itzulpen-ordain) bakan edo arraroak aipatzen dituzte, eta baztertu egiten dituzte, aldiz, garrantzi edo maiztasun handikoak. Orain, maiztasun neurketak hiztegitintzaren tresneriaren ohiko osagai bihurtu dira, eta corpus datuetan oinarritutako hitzen profilek (*word sketches*, *Wortprofile*) hitzaren agerkideei eta erlazio sintaktikoei buruzko ebidentzia ematen digute (Kilgarriff & Tugwell 2002). Maiztasun datuak, hitzaren agertokietako testuinguruak islatzen dituzten konkordantziak eta hitz-profilak segundo batean baino denbora gutxiagoan sor eta ikusaraz daitezke, dagokion hizkuntzan testu-corpus behar bezain handiak eta hizkuntzaren prozesamendurako tresnak eskuragarri izanez gero.

Hiztegian erabilera-adibideak emateko, aurreko garaietan eskuz batu behar izaten ziren aipuak, zegokien lemaren arabera *zapata kaxa* famatueta gorde eta hiztegia editatzean ateratzen zirenak. Corpusei esker, lipar batean sor ditzakegu hitz baten inguruko testuinguru edo konkordantziak, hiztegitintza-lanetan erabiltzeko. Baina hiztegitintzari hiztegi baten edukia editatzen laguntzeaz gain, beste zeregin bat betetzen dute corpus-datuak hiztegitintza arloan: gero eta sareko

hiztegi gehiagotan, erabiltzaileari erakusten zaion hiztegi sarrerari (*editoriala* eta *estatikoa* dei diezaioketun horri) corpus datuetan oinarritutako edukiak gehitzen zaizkio, erabiltzaileak bilaketa gauzatzen duen momentuan sortuta, hala nola corpusetatik erauzitako adibide-esaldiak, bilatutako hitza dutenak. Hiztegi elebidun batentzat prestatutako datu-sortak direnean, esaldi bakoitzak dagokion itzulpena du alboan; hau da, konkordantzia edo KWIC (*keyword in context*) elebidun paraleloa. Paragrafo- edo esaldi-bikote elebidunak itzulpen-corpus batetik erauzi ahal izateko, corpusa lerrokatu behar da, hots, zein segmentu zein segmenturen itzulpena den zehaztu. Paragrafo edo esaldi mailan lerrokatutako itzulpen-corpusari, corpus paraleloa (*parallel corpus*) deitzen diogu.

Baliabide askoko hizkuntza *handietan*, alemanez, esaterako, hiztegiak, baita beste edozein erabiltzailek ere, corpusetatik erauzitako datu-sortak ematen dituzten maila goreneko baliabideak ditu eskuragarri, DWDS hiztegi-atariaren eskaintza, kasu.<sup>1</sup> *Corpus hiztegi-gintza* batek zer eragin duen hiztegi-gintza lanetan eta corpusetako datuek zenbateraino aberastu ditzaketen sareko hiztegi baten edukiak, horra hor artikulu honetan interesatuko zaizkigun galderak.

## 2. Corpus elebakarrak hiztegi-gintza elebiduneko lanetan

Gure interesgune nagusia den hiztegi-gintza elebidunean murgildu baino lehen, ikus dezagun labur zer ekarpen zor diogun corpus hiztegi-gintzari elebakarrari dagokionez. Alde batetik, hiztegi baten lematagia eraikitzeke lanak aipatuko ditugu, eta bestetik, lema bakoitzari dagozkion entitate sintaktikoak (*kategoria gramatikalak*) definitzekoak. Bi kasuetan, corpusetatik erauzitako maiztasun-datuak dira lanean erabiltzeko ebidentzia ematen digutenak. Horretaz gain, jakina, erabilera-adibideak dira corpusek ematen dituzten datu garrantzitsuak, euskara elebakarrentzat corpusetan oinarritutako *Egungo Euskararen Hiztegiak*<sup>2</sup>, esaterako, baliatzen dituenak. Honetara hurrengo atalean itzuliko gara, hiztegi-gintza elebidunari begira.

Hitzen erabilera-maiztasuna eta hitz haiei dagozkien hiztegi-sarreraren ikustaldi-kopurua erlazionatuta daude neurri handi batean (De Schryver et al. 2010; Wolfer et al. 2014). Horrek esan nahi du jokabide zentzuduntzat jo dezakegula hiztegiaren lematagia eraikitzeke lanak corpusetan oinarritutako maiztasun-zerrenda batetik abiatzea, hiztegi berri bat sortzean, bederen: erabiltzaileak maiz agertzen diren hitz hauek bilatuko ditu hiztegiaren. Aldi berean, datuak eskura edukiz gero, le mari lotutako maiztasun-informazioa eman ahal zaio erabiltzaileari hiztegi-sarreraren bertan, ingelesezko hizkuntza-ikasleentzako hiztegi-gintzan aspalditik egiten den bezala (Kilgarriff 1997). Alemanez, esaterako, badaude era profesional batean hiztegi-gintzari lehengia eskaintzeko asmoz sortutako maiztasun-zerrendak.<sup>3</sup> Euskaraz, maiztasun hiztegiak

- 
1. <http://dwds.de/>. Atari honek, puntako baliabide eta teknologietan oinarrituta, alemanezko zenbait hiztegiaren edukia, corpus konkordantziak eta hitz-profilak eskaintzen ditu, erabiltzaileak bere interesaren arabera molda dezakeen modu batean.
  2. Ikus <http://www.ehu.es/eeh/>.
  3. Ikus <http://www1.ids-mannheim.de/kl/projekte/methoden/derewo.html>. Maiztasun-zerrenda guztiak lizentzia librekoak dira.

argitaratu dira, corpus batean maizen agertzen diren hitzen bildumak.<sup>4</sup> Guk ere proposatua dugu lehenago ere zer metodo erabil daitekeen eskuragarri ditugun euskarazko corpus handienetarik erazitako maiztasun-zerrendetan oinarritutako lemategi bat sortzeko (Lindemann & San Vicente 2015). Euskarazko hiztegitan sarrera izanik ere egungo corpusetan agertzen ez diren lema-multzoak zehaztu ditugu, baita alderantziz ere: egungo corpusetan bai, baina gaur arteko hiztegitan agertzen ez zaizkigunak. Corpusetan agertzeak hitzaren erabilera egiaztatzen duela onartzen badugu, hutsetik abiatzen den hiztegitintza-asmu batentzat garrantzi handiko ebidentziatzat jotzen dugu bietan agertzea: hiztegitan eta corpusetan. Batez ere euskarazko web-corpusak lexema eta adiera berri asko dituztela azpimarratzen du Leturiak (2014). Corpusetan bai, baina hiztegitan (oraindik) agertzen ez diren euskarazko hitz haiek, lema bihurtzeko hautagai direnak, web-corpus batek dituen datu-masa erraldoian ezkututzen dira, eta ez da erraza haiek iragaztea: eskuzko lana ezinbestekoa da zeregin honetan.

Bestalde, lema entitate sintaktiko gisa identifikatzeko ebidentzia eskaintzen digute corpusak. Izan ere, corpus bateko edukia prozesatzeko pauso bat hitz guztien etiketatze morfosintaktikoa da, *LemmaTizer-PoS-Tagger* deituriko tresna batez gauzatzen dena. Prozesatze horren funtzioa bikoitza da: alde batetik, testu-hitz bakoitzari dagokion forma kanonikoa zehaztea, hots, hiztegian sarrera-buru izan daitekeen lema zehaztea; eta bestetik, hitzaren kategoria (eta azpikategoria) gramatikala definitzea. Ikus dezagun adibide bat, euskarazko *alegia* formari web-corpus handi batean<sup>5</sup> zehar *EusTagger* etiketatzailerak (Aduriz et al. 1996, 2007) esleitu dizkion kategoriarik erakusten dituenak:

1. taula: Elh200 corpusetik erazitako maiztasun datuak *alegia* lemaentzat azpikategoria mailan

posizioa maiztasun-zerrendan	agerpen kopurua	kategoria gramatikala
618	41.106	lokailua
3.882	4.208	izen arrunta
10.407	921	izen berezia: leku izena
1.440.023	1	adjektiboa
1.880.968	1	adberbioa

Agerpen kopururako atalase minimoa finkatzen badugu,<sup>6</sup> *alegia* lema dagokion hiru hiztegi-sarrera (edo hiru ataleko sarrera bat) izan beharko ditu hiztegiak, datu horien arabera:

4. Hiru argitalpen ditugu: 1977ko corpus batean oinarritutako Sarasolaren (1982) lana, ahozko corpus batetik abiatzen den Etxebarria eta Mujikarena (1987), eta, hirugarrenik, UZEI (2004) maiztasun hiztegia, *XX. Mendeko Corpora* oinarria duena.
5. Datu hauek Elhuyar200 web-corpusetik erauzi ditugu. Corpus horren garapena Igor Leturiari (2014) zor zaio.
6. Sinclair (2005) lanean 20 agerpen proposatzen dira, hiztegitintza lanetan esanahia definitzeko ebidentzia nahikoa izateko. Etiketatzaile morfosintaktikoak *alegia* forma kasu banatan adjektibo eta adberbio gisa etiketatzen izana, beraz, ez da nahikoa *alegia* lemaentzat kategoria hauetako hiztegi-sarrerak egiteko, etiketatze hori zuzena ala okerra (prozesamendu errorea) den gorabehera.

*alegia* lokailuarena, *alegia* izen arruntarena, eta *alegia* (*Alegia*) herri izenarena. Ikusten dugunez, *EusTagger* tresnak maiztasun-datuak ematen dizkigu azpikategoria gramatikalari dagokionez (izen arrunta, izen berezia, etab.). Entitate sintaktiko bakoitzerako, beraz, maiztasun-datuak ditugu, hiztegi-sarrerara batean ere eman daitezkeenak.

### 3. Corpus paraleloak eta hiztezigintza

#### 3.1. Euskararekiko corpus paraleloak: artearen egoera

Elhuyar Fundazioaren web-corpusen atalean euskara-gaztelaniatzko corpus paraleloa dago kontsultagai.<sup>7</sup> Corpora automatikoki osatu da, eduki elebiduna duten domeinuak sarean bilatuz, eta domeinu horietatik elkarren itzulpena diren esaldiak erauziz. 18.753.613 testu-hitz ditu (7.891.104 euskaraz eta 10.862.509 gaztelaniatz). 659.630 segmentu elebidun ditu. Egungo euskara-gaztelaniako corpus paralelo handiena da, eta edonork erabil dezake bilaketak egiteko interfazea.

Euskal Herriko Unibertsitateko Euskara Zerbitzuak 2002. urte ingurutik argitaratzen dituen ikasliburuen itzulpenez osatua da *EHUskaratuak* corpus paraleloa. Giza zientzia, gizarte-zientzia, ingeniariaritz eta teknologia, osasun-zientzia eta zientzia zehatz eta materiaren zientzia arloetako itzulpenak biltzen ditu, ingelesetik, frantsesetik eta gaztelaniatik burutuak; corpus paralelo eleaniztuna da, beraz. Sareko interfaze bat erabilgarri dago UPV/EHUko webgunean.<sup>8</sup> Corpora Elhuyar Hizkuntza eta Teknologia garatu du, eta itzulpen berriak gehituko zaizkio corpusari etorkizunean.

#### 3.2. Euskara-alemanezko corpus paraleloak

Euskara-alemana hizkuntza-bikoteari bigarren mailakoa dei diezaiokegu, euskararen barne-erdarak (gaztelera eta frantsesa) edo kanpo-erdara nagusia (ingelesa) kide diren bikoteen atzean, hiztezigintza nahiz corpusgintzari begira. Atal honetan ikusiko dugunez, bigarren mailako bikote honetan kalitatezko corpus paraleloak sortzeko ahaleginak egin berri dira. Bigarren maila honetako beste bikote guztietan, berriz, egiteke geratzen zaizkigu oraindik, guk dakigunez.

EHUko Letren Fakultatean kokatutako TRALIMA/ITZULIK ikerketa taldean, alemanezko literatura-lanak nahiz euskarazko itzulpenak biltzen dituen corpus paralelo bat sortu du Uribarri, Sanz eta Zubillaga hirukoteak (Sanz et al. 2015). Corpusak itzulpen zuzenak eta zeharkakoak biltzen ditu. Alemanetik zuzenean euskaratu diren testuak alemanezko jatorrizkoan nahiz euskarazko bertsioan daude jasota; zeharka burutu diren itzulpenak, aldiz, jatorrizkoan, gaztelerazko *zubi-bertsioan* eta euskarazkoan. Azken kasu hauetan, beraz, corpora hirueleduna

7. Ikus <http://webcorpusak.elhuyar.org/cgi-bin/kontsulta2.py>.

8. Ikus <http://ehuskaratuak.ehu.eus/bilaketa/>.

da. Corpora osatzen duten testuetatik, alemanezko 81 testu literario eta euskarazko itzulpen zuzenak izan ditugu eskuragarri, paragrafo edo esaldi mailan lerrokaturik, hiztegitintza arloko gure esperimenduetan erabiltzeko. Euskarazko testuen lematizatzea *EusTagger* tresnaren bidez gauzatu dugu<sup>9</sup>. *SketchEngine* (Kilgarriff et al. 2004) tresnaren bitartez, alemanezko testuen lematizatzea (*TreeTagger*, (Schmid 1995)) burutu eta, testu paralelo guztiak aplikazioan kargaturik, bilaketak egin eta konkordantzia paraleloak sortu ditugu. Horrela sortutako corpus paraleloak bina milioi testu-hitz inguru ditu euskaraz eta alemanez. Bestetik, Bibliaren alemanezko eta euskarazko bertsio bana erabili dugu corpus paraleloa sortzeko.<sup>10</sup> Bibliatik corpus paraleloak sortzea nahiko ariketa erraza da, bertset-kodeen bitartez lerroka baitaiteke. Bibliaren itzulpena hainbat eta hainbat hizkuntzatan, eta, horretaz gain, hizkuntza bakoitzean maiz bertsio batean baino gehiagotan eskuragarri izatea oinarri egokia da era askotako konparaketak egiteko (Resnik et al. 1999; Nida & Taber 2003). Corpusak 700.000 bat testu-hitz ditu hizkuntza bakoitzean.

Eskuz hautatutako iturri literario hauetaz gain, euskara-alemanezko corpus paraleloak modu eraginkor batean osatzeko balioko luketen beste zenbait datu-iturri ere aipa daitezke: ikus-entzunezkoen azpтитuluak eta software lokalizazioak.<sup>11</sup>

### 3.3. Corpus paraleloak hiztegitintza lanetan

*EuDeLex* izeneko euskara-alemanezko gure egitasmo berriaren garapen prozesuan, aurreko atalean deskribatutako corpus elebiduna erabiltzen dugu. Batetik, konkordantzia paraleloak baliatzen ditugu alemanezko lemaren euskarazko ordainak definitzeko lanetan. Gure esperientziaren arabera, batez ere kontzeptu abstraktuen kasuetan lagungarriak izan dira itultzaileen proposamenak euskarazko ordain egokienak hiztegi-sarrera elebidunerako hautatzeko, baita dagokion lema duten esapide eta lokuzioak itultzeko ere. Bestetik, euskara-alemanezko glosarioak zirriboratu ditugu datu haietan oinarriturik: esaldi- edo paragrafo-mailako lerrokaketatik hitz-mailako lerrokaketara heltzen diren tresna automatikoen bitartez, aleman-euskarazko itzulpen-ordainen bikote izateko hautagaiak definitu ditugu (ikus Lindemann et al. 2014).

### 3.4. Corpus paraleloetatik erauzitako adibide-perpauak hiztegi-sarreraren osagai

Corpus paraleloetatik erauzitako konkordantziak hiztegi elebidunetako sarreretan islatzeko aukera 1990eko hamarkadaren hasieratik aipatu izan den arren (lehenengo aipamenak Atkins 1996; Dickens & Salkie 1996), aukera berri hori gauzatu izanaren zenbait adibide baino ez

---

9. EHUko IXA taldeko Gorka Labakari zor dizkiogu euskarazko lematizatze eta etiketatze lanak.

10. Elhuyar Fundazioko Xabier Saralegiri zor zaio lan hori, ikus Lindemann et al. (2014).

11. OPUS egitasmoak euskarazko eta beste hainbat hizkuntzako datu paraleloak biltzen ditu, iturri haietatik batik bat (Lindemann 2012). Euskara-alemanezko konbinazioan, adibidez, OPUS datu-basean batutako testu paraleloa aintzat hartzeko tamaina batera irits liteke datozen urteotan.

ditugu egun arte. Gorago esan bezala, corpusak hitzen adiera nagusiak edukitzen ditu, hiztegi sarrera editorial batean maiztasun gutxiagoko adieren artean *ezkutaturik* gera litezkeenak. Bes-tetik, hitz baten adiera helburu hizkuntzan islatzeko, itzultzaileek maiztasun gutxiko ordainak hautatzen dituzte maiz, hiztegi-sarrera editorial batean agertzen ez direnak. Wolfgang Teubert autorearen hitzetan, “the translator’s design space is much larger than the language-neutral conceptual ontology (or the traditional bilingual dictionary) would leave us to believe” (Teubert 2002). Gainera, itzultzaile bat ez da mugatzen hitz bat hitz batez ordezkatzera; aitzitik, egitura sintaktikoa ere moldatu, laburbildu, garatu, edo helburu hizkuntzan egokia gertatzen den esapide batez itzul dezake, esaldi osoa modu egoki batean islatzeko bere arduraren arabera. Giza itzultzaileak baino sor ez ditzakeen moldapen horiek dira kalitatezko corpus paralelo ba-tek eskaini dituen altxorak.

2. taula. Corpus paralelotik erauzitako alemanezko *Ärger* lema-aren konkordantziak

<i>Konkordantzia DE</i>	<i>Konkordantzia EU</i>	<i>Transformazio sintaktikoa</i>
Aber auf dieser Route <b>hatte</b> ich noch nie <b>Ärger</b> mit ihnen!	Baina bide honetan ez <b>dut</b> inoiz haiekin <b>arazorik izan!</b>	V+N > V+N
Nun war der <b>Ärger</b> in Capricorns Stimme nicht mehr zu überhören	Dagoeneko nabaria zen Kaprikornioaren ahotsean <b>haserrea</b> .	N > N
Meggie konnte nicht verhindern, dass ihre Stimme <b>vor Ärger</b> zitterte.	Meggiek ezin ekidin zezakeen ahotsak <b>amorruez</b> dar-dar egin zieziaion.	prep.+N > N+postp.
Meggie sah, wie sich Bastas Schultern <b>vor Ärger</b> spannten.	<b>Haserrearen haserrez</b> , Bastaren sorbaldak tenkatu egiten zirela ikusi zuen Meggiek.	prep.+N > (esapidea) adb.
Mensch, biste nicht froh, daß du den ganzen <b>Ärger</b> los bist mit den Weibern?	Eta hi zer, ez al hago pozik, andreekin hituen <b>saltsa</b> guztiak bukatuta?	N > N
Bastas Lippen wurden schmal <b>vor Ärger</b> , doch er verkniff sich eine Antwort und...	Bastak ezpainak estutu zituen <b>haserre</b> , baina erantzuteko gogoari eutsi zion.	prep.+N > adb.

Egokitzapen ezberdinak ikusten ditugu goiko taulako euskara-alemanezko adibideetan. Corpus paraleloetako datuetan, beraz, itzultzaileek jorratutako bideak azaltzen zaizkigu, itzul-tzaileak zer prozeduraz baliatu diren, kasuz kasu, testuingurua aintzat harturik; eta horixe da hiztegi-sarrera elebidun batek bere estatusunean orain arte nekez aurreikus eta isla zezakeena. Horrelako datuetan hizkuntzalari konputazionalek bilatzen dute itzulpen-makinek oraindik imitatzen ez dakitena, taulan islatu duguna, “syntactically motivated transformation rules that

explain human translation data” (Galley et al. 2004). Testuingurua duten itzulpenez, corpus paraleloz elikatzen dituzte estatistika hutsa baliatzen besterik ez dakiten itzulpen-tresna automatikoak, haiei giza itzultzaileen jokabideak *ulertarazteko* ahaleginean.

Hiztegi-sarrera editorial batean, lemak entitate sintaktiko gisa dituen izaerak (hau da, haren kategoria gramatikalak) mantendu ohi dira hizkuntza batetik bestera, helburu hizkuntzan posible den heinean, bederen.<sup>12</sup> Beraz, hiztegi-sarrera elebidunarekin batera benetako adibideen ondoan benetako itzulpenak ematea lagungarria da benetan erabiltzen den hizkuntza islatzeko, eta corpusetako konkordantzia paraleloak horretarako baliatzea jokabide egokia izan daiteke. Zalantzarik gabe, corpus paraleloaren edukien kalitatea aldagai garrantzitsua da; orain arteko gure esperientziaren arabera, liburuaren itzulpenak dira helburu horretarako emankorrenak.

### 3.5. Konkordantzia paraleloak hiztegi elektronikoetan: artearen egoera

Itzulpen-corpusetatik erauzitako datuak hiztegi-erabiltzaileari zuzenean erakustea eta horrela erabiltzailearen galderari benetako erabilera-adibide eta haien itzulpenez erantzutea joera berria dugu hiztegitintza elebidun elektronikoan. Bide berriak jorratzen dituzten bi hiztegi-atari aipatuko ditugu.

*Linguee* izeneko egitasmoa 2007an sortu zen, eta 2009tik aurrera, erabiltzaile-interfazea erabilgarri dago edonorentzat.<sup>13</sup> 2015ean, 25 hizkuntzako datuak biltzen ditu eta haien arteko oinarrizko hiztegi-sarrera elebidunak eta itzulpen-adibideak eskaintzen ditu. Itzulitako esaldien datu-baseak Interneten aurkitutako testu paraleloak ditu. Ikus dezagun adibide bat, gaztelera-ko *falda* lema polisemikoaren ingelesezko itzulpenak biltzen dituen (1. irudia):

Ikusten denez, galderaren erantzunak biltzen dituen orri honetako goiko partean, hiztegi-sarrera elebiduna dator, gaztelera-ko hitzaren ezaugarri morfosintaktikoak eta hiru adierentzako ingelesezko ordain bana ematen duena. Horretaz gain, hitzak entzuteko aukera ematen duten botoiak, eta *falda* parte duten unitate lexiko hitz anitzekoen ingelesezko ordainak ematen dira *adibide* gisa. Azpian, ingelesezko nahiz gaztelera-ko *wikipediak* ematen dituzten definizio laburrak ikus daitezke, dagozkien geztetan klikatuz. Honaino, beraz, oso eskaintza oinarrizkoa da ikusten duguna.

Hortik behera, itzulitako adibide-esaldiak datoz, eta irudi honetan haien parte txiki bat baino ez dugu islatu. Ikusten denez, gaztelera-ko hitza nahiz ingelesezko ordaina nabarmentzen dira, bilatutako hitza eta ordainak testuinguruetan zehar bizkorrago ikusteko. Honelako esaldibildumen balioa bistakoa da: hiztegi-erabiltzaileak bere kasuari hurbila zaion testuingurua bi-

12. Batzuetan, helburu-hizkuntzan posible dena azken muturreraino erabiltzen da, kategoria gramatikala mantentzearen. Horrela, *Euskal WordNet* datu-base lexikalean ingelesezko hainbat izenetako ordainak eratorritako formak dira: adibidez, hainbat *-(z)ea* eta *-tasun* izen agertzen dira euskarazko ordain gisa, ingelesezko iturria izena denean. Horietako asko ez dira eratorritako forma horretan agertzen euskarazko hiztegi-lemategietan: *connivance*, *secret approval*, *tacit consent*: ‘ados jartze’, *spraying*, *crop-dusting*: ‘aerosolez bustitze’, *randomness*, *stochasticity*: ‘aleatoriotasun’.

13. Ikus <http://www.linguee.com>.

falda con vuelo *f* — fit and flare skirt *s*

c-3

▶	Wikipedia
▼	Fuentes externas
Godet: pliegues de tela cortados a capa que se añaden en la parte baja de la falda para darle más vuelo y más peso. <span style="float: right;">🌐 <a href="http://esflamenco.com">esflamenco.com</a></span>	Godet: folds of cloth cut i to give it more fullness ai
La segunda figura lleva falda y sostiene algo que podría ser una bolsa de red. <span style="float: right;">🌐 <a href="http://paracas.se">paracas.se</a></span>	The other figure is wearin
Incluso antes de que los bebés puedan sentarse por sí mismos, pueden acurrucarse en su falda y mirar imágenes de libros con dibujos en blanco y negro o fotos de rostros. <span style="float: right;">🌐 <a href="http://es.4children.org">es.4children.org</a></span>	Even before babies can and look at board books
Corte y deshuesado: separar toda la falda del cuarto trasero desde el corte recto a la altura de la octava costilla, cortando [...]. <span style="float: right;">🌐 <a href="http://eur-lex.europa.eu">eur-lex.europa.eu</a></span>	Cutting and boning: rem hindquarter by a cut from
Si, tenemos un pantalón corto o una falda en beis, una camiseta roja y una camisa azul. <span style="float: right;">🌐 <a href="http://filmcap.org">filmcap.org</a></span>	Yes, a beige trouser or si
La mujer lleva una blusa blanca y una falda negra. <span style="float: right;">🌐 <a href="http://resources.rosefastore.com">resources.rosefastore.com</a></span>	The woman is wearing a
La villa está enclavada junto al río Matarraña, en las proximidades del límite de provincia de Tarragona, en la falda de la Reserva Natural de los Puertos de Beceite. <span style="float: right;">🌐 <a href="http://expozaragoza2008.es">expozaragoza2008.es</a></span>	The town is right next to l province of Tarragona, o Beceite.
Asimismo, las autoridades han prohibido a las mujeres usar bicicletas (un vehículo esencial para poder comerciar) y las han obligado a llevar falda. <span style="float: right;">🌐 <a href="http://daccess-ods.un.org">daccess-ods.un.org</a></span>	The authorities have als key vehicle for access to
El pistón de falda corta se mueve en un cilindro recubierto de un compuesto cerámico que reduce las pérdidas friccionales, [...]. <span style="float: right;">🌐 <a href="http://yamaha-motor-europe.com">yamaha-motor-europe.com</a></span>	The short skirt piston run reduces frictional losses,
Alojamiento en suites con piscina o pabellones de caprichosa arquitectura en la falda de cuatro volcanes humeantes. <span style="float: right;">🌐 <a href="http://gulabomejordelmundo.com">gulabomejordelmundo.com</a></span>	Accommodation in suite capricious design set on
En el armario, todavía están la blusa y la falda del primer día, nada más. <span style="float: right;">🌐 <a href="http://basque literature.com">basque literature.com</a></span>	The jacket and skirt from wardrobe.

1. irudia: *Linguee* hiztegiaren lagina



latuko du, edota kasuz kasu nola itzuli den ikusi. Jakin nahi duena, beharbada, oso galdera berezitua da, termino gisako erabilera bat, ibilgailu-motorreko *pistón de falda corta* ingelesez *short skirt piston* ote den, adibidez, eta datu-base terminologiko egokia eduki ezean, hemen bai aurkitzen du erantzun bat, kasualitatea dei badiezaiokegu ere. Gaztelera-ingelera konbinazioan badira banku terminologikoak kasu honetarako<sup>14</sup>; baliabide gutxiagoko hizkuntza-bikote baterako, aldiz, corpus paraleloak baliatzen dituen tresna bat edukitzea bereziki lagungarria dela iruditzen zaigu.

1.: ya que	2.: sin embargo	3.: por lo tanto	4.: además	5.: por otro lado
6.: aunque	7.: es decir	8.: a partir de	9.: llevar a cabo	10.: por lo que
11.: asimismo	12.: no obstante	13.: por otra parte	14.: a continuación	15.: aprovechar
16.: actualmente	17.: a pesar de	18.: so	19.: así como	20.: mediante
21.: en cuanto a	22.: realizar	23.: aportar	24.: sobre todo	25.: cumplir
26.: a su vez	27.: debido a	28.: en este sentido	29.: por su parte	30.: atentamente
31.: presupuesto	32.: conseguir	33.: seguimiento	34.: destacar	35.: siempre y cuando
36.: tener en cuenta	37.: de esta manera	38.: vigente	39.: por tanto	40.: posteriormente
41.: con respecto a	42.: facilitar	43.: funcionamiento	44.: empresa	45.: convocatoria
46.: sino	47.: si bien	48.: propuesta	49.: compromiso	50.: resumen

*Linguee* gunean bertan ematen diren ikustaldi handiko sarreren zerrenden arabera (ikus goiko zerrenda), ondorengoa bezalakoak dira oso maiz baliabide honen bitartez argitzen saiatzen diren kasuak: hizkuntza batek berezko dituen hitz anitzeko esapideak, gaztelerazko *en cuanto a* bezalakoak. Ikus ditzagun haren itzulpenak:

<p><b>En cuanto a la composición demográfica del país, el proceso de envejecimiento de la población impone un reto a la política social.</b>  <small>↳ <a href="http://daccess-ods.un.org">daccess-ods.un.org</a></small></p>	<p><b>With regard to the country's demographic composition, the ageing of the population poses a challenge to social policy.</b>  <small>↳ <a href="http://daccess-ods.un.org">daccess-ods.un.org</a></small></p>
<p><b>En esta labor se debe tener en cuenta que persisten las diferencias entre los países en cuanto a su capacidad para la difusión de información.</b>  <small>↳ <a href="http://daccess-ods.un.org">daccess-ods.un.org</a></small></p>	<p><b>Those initiatives must take into account the differences that persisted among countries' capacities to disseminate information.</b>  <small>↳ <a href="http://daccess-ods.un.org">daccess-ods.un.org</a></small></p>
<p><b>En cuanto a la cuestión de la energía, por desgracia soy bastante pesimista.</b>  <small>↳ <a href="http://europarl.europa.eu">europarl.europa.eu</a></small></p>	<p><b>When it comes to the energy issue, I am unfortunately rather pessimistic.</b>  <small>↳ <a href="http://europarl.europa.eu">europarl.europa.eu</a></small></p>
<p><b>Dado que la búsqueda se inicia en cuanto empieza a escribir, si por ejemplo escribe "B", todos los archivos con nombres [...]</b>  <small>↳ <a href="http://windows.microsoft.com">windows.microsoft.com</a></small></p>	<p><b>The search begins as soon as you begin typing-so if you type "B," for example, all the files with names starting with the [...]</b>  <small>↳ <a href="http://windows.microsoft.com">windows.microsoft.com</a></small></p>
<p><b>En la actualidad, ellos deciden en cuanto a las ofertas sobre la base de la calidad y el precio.</b>  <small>↳ <a href="http://gopacnetwork.org">gopacnetwork.org</a></small></p>	<p><b>At present, they decide on bid offers using the bases of quality and price.</b>  <small>↳ <a href="http://gopacnetwork.org">gopacnetwork.org</a></small></p>
<p><b>[...] Unidos se encuentra mucho más adelantado que el resto de las naciones en cuanto al gasto en servicios de salud per cápita, ocupa el lugar 48 en expectativa de vida, dijo.</b>  <small>↳ <a href="http://forumfed.org">forumfed.org</a></small></p>	<p><b>At the same time, the United States remains far in front of all other nations in per capita health care spending while it ranks 48th in life expectancy, she said.</b>  <small>↳ <a href="http://forumfed.org">forumfed.org</a></small></p>
<p><b>En cuanto a mí, no puedo pensar en hacer películas en las que no me involucre, prefiero esperar un proyecto que me estimule de verdad.</b>  <small>↳ <a href="http://cineuropa.mobi">cineuropa.mobi</a></small></p>	<p><b>It is impossible for me to get involved in films that I don't like so I just wait for a project that really tickles my fancy.</b>  <small>↳ <a href="http://cineuropa.mobi">cineuropa.mobi</a></small></p>
<p><b>[...] diversifica la oferta de educación se observan desigualdades de una magnitud sin precedentes en cuanto al acceso y la calidad.</b>  <small>↳ <a href="http://unesdoc.unesco.org">unesdoc.unesco.org</a></small></p>	<p><b>But as educational demand increases and supply diversifies, disparities of unprecedented proportions can be observed in respect of access and quality.</b>  <small>↳ <a href="http://unesdoc.unesco.org">unesdoc.unesco.org</a></small></p>
<p><b>El parlamento es útil a la burguesía en cuanto logra crear la ilusión de que los obreros eligen a quienes los malgobierman.</b>  <small>↳ <a href="http://ibrp.org">ibrp.org</a></small></p>	<p><b>Parliament is useful to the bourgeoisie in that it gives the illusion that workers choose who is to misrule them.</b>  <small>↳ <a href="http://ibrp.org">ibrp.org</a></small></p>

14. Egun, Europar Batasuneko hizkuntza ofizialak batzen dituen IATE datu-base terminologikoa puntako baliabidea dugu; ikus <http://iate.europa.eu>.

Ikusten eta espero dugunez, *en cuanto a* esapideak askotariko ordainak ditu ingelesez, testuinguruaren arabera. Ordain bakar bat ez da errepikatzen. Kasu batzuetan, ohiko hiztegi-sarrerera batean aipatuko litzatekeen ordain zuzenik ez da, bigarren adibidean duguna, kasu: genitibozko egitura baten bitartez itzultzen da gaztelerazko esapidea. Aipatzekoa da, halaber, gaztelerazko laugarren adibideko *en cuanto* egiturak, grafiaren hurbiltasuna gorabehera, ez duela zerikusi semantikorik bilatutako *en cuanto a* esapidearekin; horra hor atari honetan erabilitako metodoaren muga bat. Erabiltzailearentzako onura, berriz, agerian dago: ohiko hiztegi-sarrerera batean agertuko ez litzazkiokeen itzulpenak, testuinguru ezberdinak dituzten itzulpenak datozkio hemen.

Bestetik, maizen gauzatzen diren bilaketa hauek ikusita, ondorengo susmoa sortzen zaigu: hitz anitzeko unitate ugari dira zerrenda hauetan; erabiltzaileek hitz anitzeko unitate haiek beste hiztegi batzuetan, testuinguruaren duten itzulpenak ematen ez dituzten hiztegiengan, aurkitzea espero ez dutelako ote? Hiztegien erabileraren inguruko ikerketa-galdera irekia dugu hemen. *Linguee*, behintzat “a translators’ favourite” bilakatu da epe motz batean (Kilgarriff 2013:92).

*Linguee* atariaren garatzaileek haien metodoen berri ematen ez badute ere, bistan da logaritmoak ere erabiltzen direla corpusak dituen adibide-esaldi guztien artean kontsulta bati erantzuteko egokienak hautatzeko, ez baitago esaldi motzegirik edo luzezegirik ezta oso testuinguru edo hitz bitxia daukan esaldirik. Adibide *onak* hautatzeko logaritmo automatikoak garatu izan dira, eta zenbait kasutan publikoak dira haiek erabiltzeko arautegiak, Kilgarriff (2008) adibide.

*Glosbe* hiztegi-ataria<sup>15</sup> euskarazko esaldiak eta hiru erdara nagusiez gainerako hainbat hizkuntzarako itzulpenak eskaintzen dituen baliabide bakarra da egun, guk dakigunez. Izan ere, hizkuntza-bikote askorentzako itzulpenak ditu atari honetako datu-baseak. Euskaratik edo euskarara ehun bat hizkuntzatarara ditu loturak, eta arlo jakin batzuetarako baliabide baliotsutzat jo dezakegu dagoeneko, informatikarekin zerikusia dutenetarako, bereziki. Izan ere, software lokalizazioetako edukietan oinarritzen dira *Glosbe* datu-basearen euskarazko laginak batez ere.

#### 4. Ondorioak

Itzulpen-corpusak hiztegi-gintzan baliatzea aberasgarria da, zalantzarik gabe, hiztegi-gilearentzat nahiz hiztegi erabiltzailearentzat: corpusak tamaina eta osakeraren arabera eskaini ditzakeen adibide eta itzulpenek ohiko hiztegi elebidunen egitura aberasten dute. Adibide-esaldiak eskuz hautatzea eta dagokien hiztegi-sarreran txertatzea, corpusen garaiaren aurretik egiten zen legez, lan erraldoia izaten zen. Gaur, corpus-konkordantziek errazten diote hiztegi-gileari eduki egokienak hautatzeko lana. Ibon Sarasolaren *Egungo Euskararen Hiztegia* corpus elebakar batek eskuz hautatutako erabilera-adibideak baliatzearen etsenplu bikaina da. Adibide egokienak eskuz aukeratzeko baliabiderik ez dagoenean edo ahalik eta adibide gehien eman nahi direnean, erabiltzaileak bilaketa egiten duen momentuan sorturiko corpus-konkordantziak erants

15. Ikus <http://eu.glosbe.com/eu/>.

dakizkioke hiztegi-sarrera estatiko bati, goian aipatutako DWDS atarian gertatzen den legez. Corpus paraleloetatik automatikoki erauzitako itzulpen-adibideak hiztegi elektronikoko baten sarreretan zuzenean txertatzea joera berria da; eskuz landu gabeko hizkuntza-bikote askotan, corpus paraleloetan oinarriturik automatikoki sortutako hiztegiak dira eskuragarri dauden baliabide bakarrak. Dena den, kalitatezko corpus hiztegegintza batean ezin dugu ahaztu eskuzko lanaren premia, metodo automatikoen bitartez sortutako emaitza desegokiak iragazteko eta hutsuneak atzeman eta betetzeko orduan, batik bat.

## 5. Aipamenak

- ADURIZ, I., ALDEZABAL, I., ALEGRIA, I., ARTOLA, X., EZEIZA, N. & URIZAR, R. (1996). EUSLEM: A lemmatiser/tagger for Basque. In *Proceedings of EURALEX 1996*. Göteborg: Göteborg University, pp. 17–26.
- ADURIZ, I., ALEGRIA, I., ARRIOLA, J., ARTOLA, X., DÍAZ DE ILARRAZA, A., EZEIZA, N., GOJENOLA, K. & MARITXALAR, M. (2007). Different Issues In The Design Of A Lemmatizer/Tagger For Basque.
- ATKINS, B.T.S. (1996). Bilingual dictionaries: Past, present and future. In *Proceedings of EURALEX 1996*. Göteborg: Göteborg University, pp. 515–546.
- ATKINS, B.T.S. & RUNDELL, M. (2008). *The Oxford Guide to Practical Lexicography*. Oxford University Press.
- DICKENS, A. & SALKIE, R. (1996). Comparing Bilingual Dictionaries with a Parallel Corpus. In *Euralex'96 Proceedings*. Göteborg: Göteborg University, pp. 551–559.
- ETXEBARRIA, J.M. & MUJIKA, J.A. (1987). *Euskararen oinarritzako hiztegia : maiztasun eta prestasun azterketa*. Gasteiz: Eusko Jaurlaritzaren Argitalpen Zerbitzu Nagusia.
- GALLEY, M., HOPKINS, M., KNIGHT, K. & MARCU, D. (2004). What's in a translation rule. In *Proceedings of HLT/NAACL*. Boston, pp. 273–280.
- KILGARRIFF, A. (1997). Putting frequencies in the dictionary. In *International Journal of Lexicography*, 10, 135–155.
- KILGARRIFF, A., HUSAK, M., MCADAM, K., RUNDELL, M. & RYCHLÝ, P. (2008). GDEX: Automatically finding good dictionary examples in a corpus. In *Proceedings of EURALEX 2008*.
- KILGARRIFF, A., RYCHLY, P., SMRZ, P. & TUGWELL, D. (2004). The Sketch Engine. In *Proceedings of EURALEX 2004*. Lorient, France, pp. 105–116.
- KILGARRIFF, A. & TUGWELL, D. (2002). Sketching words. In Corréard, M.-H. (ed.) *Lexicography and natural language processing: a festschrift in honour of BTS Atkins*. Euralex, pp. 125–137.

- KLOSA, A. (2007). Korpusgestützte Lexikographie: besser, schneller, umfangreicher. In Kallmeyer, W. & Zifonun, G. (eds.) Sprachkorpora. Datenmengen und Erkenntnisfortschritt, Jahrbuch des Institut für Deutsche Sprache. Walter de Gruyter, pp. 105–122.
- KRISHNAMURTHY, R. (2002). The Corpus Revolution in EFL Dictionaries. In *Kernerman Dictionary News*, 10.
- LETURIA, I. (2014). *The Web as a Corpus of Basque* (PhD Thesis, UPV/EHU).
- LINDEMANN, D. (2013). Bilingual Lexicography and Corpus Methods. The Example of German-Basque as Language Pair. In *Procedia - Social and Behavioral Sciences*, 95, 249–257.
- LINDEMANN, D. & San Vicente, I. (2015). Euskarazko maiztasun lematagia gaurko teknologien ikuspuntutik. In Fernández Fernández, B. & Salaburu Etxeberria, P. (eds.) Ibon Sarasola, gorazarre. Homenatge, homenaje. Bilbo: UPV/EHU, pp. 441–456.
- LINDEMANN, D., Saralegi, X., San Vicente, I., Manterola, I. & Nazar, R. (2014). Bilingual Dictionary Drafting. The example of German-Basque, a medium-density language pair. In *Proceedings of the XVI Euralex International Congress*. EURALEX, Bolzano: EURAC, pp. 563–576.
- MITXELENA, K. & SARASOLA, I. (1988). *Diccionario general vasco - Orotariko euskal hiztegia*. Euskaltzaindia; Editorial Desclée de Brouwer.
- NIDA, E.A. & TABER, C.R. (2003). *The Theory and Practice of Translation*. 4th ed. Brill.
- RESNIK, P., OLSEN, M.B. & Diab, M. (1999). The Bible as a parallel corpus: Annotating the “Book of 2000 Tongues.” In *Computers and the Humanities*, 33, 129–153.
- RUNDELL, M. & STOCK, P. (1992). The corpus revolution. In *English Today*, 8, 45–51.
- SANZ, Z., ZUBILLAGA, N. & URIBARRI, I. (2015). Estudio basado en corpus de las traducciones del alemán al vasco. In Sánchez Nieto, M.T. (ed.) *Corpus-based Translation and Interpreting Studies: from description to application*. Berlin: Frank & Timme, pp. 211–235.
- SARASOLA, I. (1982). *Gaurko euskara idatziaren maiztasun-hiztegia: 1977ko corpus batean oinarritua*. Donostia: Gipuzkoako Aurrezki Kutxa Probintziala.
- SCHMID, H. (1995). Improvements In Part-of-Speech Tagging With an Application To German. In *Proceedings of the ACL SIGDAT-Workshop*. Dublin, pp. 47–50.
- DE SCHRYVER, G.-M., JOFFE, D., JOFFE, P. & HILLEWAERT, S. (2010). Do dictionary users really look up frequent words?—on the overestimation of the value of corpus-based lexicography. In *Lexikos*, 16.
- SINCLAIR, J. (2005). Corpus and text-basic principles. In Wynne, M. (ed.) *Developing linguistic corpora: A guide to good practice*. Oxford: Oxbow Books, pp. 1–16.

- SVENSÉN, B. (2009). *A handbook of lexicography: the theory and practice of dictionary-making*. Cambridge: Cambridge University Press.
- TEUBERT, W. (2002). The role of parallel corpora in translation and multilingual lexicography. In Altenberg, B. & Granger, S. (eds.) *Lexis in Contrast: Corpus-Based Approaches*. John Benjamins Publishing, pp. 189–214.
- TIEDEMANN, J. (2012). Parallel Data, Tools and Interfaces in OPUS. In *Proceedings of the Eighth conference on International Language Resources and Evaluation*. LREC 2012, Istanbul, pp. 2214–2218.
- UZEL. (2004). *Maiztasun Hiztegia*. Donostia: UZEL.
- WOLFER, S., KOPLINIG, A., MEYER, P. & MÜLLER-SPITZER, C. (2014). Dictionary Users do Look up Frequent and Socially Relevant Words. Two Log File Analyses. In *Proceedings of the XVI Euralex International Congress*. Bolzano/Bozen: Eurac, pp. 281–290.

## Resumen

La lingüística de corpus ha revolucionado la lexicografía: los datos de frecuencia, los ejemplos de uso y las traducciones son, entre otros, los recursos que nos ofrecen los extensos corpus de texto. La frecuencia de uso de una palabra y la frecuencia de visitas de diccionario de la palabra en cuestión están en relación; además, estos datos se pueden integrar en la entrada misma de diccionario. Por lo tanto, es interesante tomar listas de frecuencia como punto de partida para definir un lecionario de diccionario.

En la lexicografía bilingüe se pueden utilizar los datos extraídos de corpus paralelos. En los últimos años se han creado varios corpus paralelos con euskera, como por ejemplo el corpus literario alemán-euskera. De estos corpus paralelos podemos extraer ejemplos de uso y sus respectivas traducciones, los que pueden resultar de ayuda en el trabajo lexicográfico, así como de información para el usuario. Esa última opción se está difundiendo en la actualidad a través de algunos portales de diccionario.

## Résumé

La linguistique de corpus a constitué une véritable révolution dans le monde de la lexicographie, en effet, les larges corpus de texte offrent de multiples ressources telles que les données de fréquence, les exemples d'utilisation et les traductions. La fréquence d'utilisation d'un terme et la fréquence de visites de dictionnaire du terme en question sont en étroite relation ; en outre, ces données peuvent être intégrées dans l'entrée même de dictionnaire. Il est, par conséquent, très intéressant de prendre les listes de fréquence comme point de départ pour définir une liste d'entrées de dictionnaire.

Dans la lexicographie bilingue on peut utiliser des données extraites de corpus parallèles. Ces dernières années plusieurs corpus parallèles comprenant l'euskara ont été créés, comme par exemple le corpus littéraire allemand-euskara. Nous pouvons extraire de ces corpus parallèles des exemples d'utilisation et leurs traductions respectives, qui peuvent être très utiles dans un travail lexicographique, tout en apportant une information précieuse à l'utilisateur. Actuellement, cette dernière possibilité est en train de se répandre grâce à certains sites de dictionnaire.

## Abstract

Corpus linguistics have revolutionized lexicography: data on frequency, usage examples, and translations are, among other things, the resources offered by extensive text corpora. The frequency of use of a word and the frequency of dictionary visits to that word are related; furthermore, these data can be integrated into the dictionary entry. Therefore, it is interesting to take frequency lists as a point of departure in defining the entry list of a dictionary.

Data obtained from parallel corpora can be used in bilingual lexicography. In recent years, a number of parallel corpora have been created with Basque, such as the literary German-Basque corpus. From these parallel corpora we can extract examples of use and their respective translations, which can be helpful not only in lexicographic work, but also in providing information to the user, a service that is becoming more and more common through the use of dictionary websites.

