

# Hizkuntzaren tratamendu automatikoa: helburuak eta abiaburuak

Inaki Alegria, Xabier Artoia eta Kepa Sarasola<sup>1</sup>

## Sarrera: lengoia naturalaren prozesamenduren helburuak

Inpentaren sorkuntzak hizkuntzaren tratamendua eta zabalgunza irauli zuen moduan, mende honean sortu den ordenadoreak horren pareko iraultza ekarri du. Testu-prozesaketarako baliabide berria den ordenadoreak erraztasun handiak eskaintzen ditu gaur egun testuak kopiatu, osatu eta zuzentzeko, eta baita mila formatu edo itxura desberdinetan aurkeztu ahal izateko ere. Baina testu-edizioko tresna horiek baino askoz ere laguntza hobekia dira merkatumen eta are laguntza bereziagoak bilatzen dira ikertokietan. Ordenadorearen bitartez hizkuntzaren tratamendua egiten duen aplikazioak eta programak gero eta ugariago dira. ordenadoreko komunikazioa egunero erabiltzen dugun hizkuntzaren bitartez egin ahal izatea gero eta normalago izango baita. Beste alde batetik, gizarte eleantizak hizkuntza diferentean artean egin behar izaten dituen joan-etorriak leuntzeko ere aparteko lagun izango dugu ordenadorea. Gainera, telekomunikazioetan gertatutako aurrerapen izugarriak eragin duen Internet fenomenoa areagotu egin du hizkuntzaren tratamendu automatikoren beharra; zeren eta nahiz eta sarearen bidez informazio kopuru izugarria lortu ahal izan, ez baita erraza bilatzen dugun informazioa aurkitzea, eta informazioa ondo selekzionatzeko tratamendu linguistikoa lagungarria baino ezinbestekoa da.

Hizkuntzaren tratamendu automatikoren inguruko ikertarlanari *Lengoia Naturalaren Prozesamendua* (LNP) esaten diogu, nahiz eta batzuetan, hizkuntzalarietako ikuspuntua garrantzitsua denean batez ere, *Linguistika Komputazionala* ere esan. Hizkuntzaren industria oso bat sortzen ari da, ordenadoreaz baliatuz hizkuntza tratatzea helburu duena. Hizkuntzaren teknologiaz hitz egiten da dagoeneko. Teknologia horren oinarrian ikerkuntza dago, hizkuntzaren tratamendu automatikoren arloko ikerkuntza, alegia. Horiek guztiak dira artikulu honen aztergaiak; hasieratik argi utzi behar dugu, ordea, ez garaela arloko euskarazko softwareaz, hau da, euskaraz erabili daitezkeen ordenadore-programei buruz<sup>2</sup>, ezta orokorrean euskarrik informatikaren munduan egun duen leku eskasaz ere<sup>3</sup>.

Dena dela, azken urteorako lorpen hauek mugatuta dagozute. Bost urteko edozein urte hitz egiten eta ulertzen ondo moldatzen denez, hizkuntza erabiltzea lan erraza dela pentsatzen dugu, baina hori ez da horrela. Lengoia sortzea eta ulertzea oso prozesu konplexuak dira eta gaur egungo

<sup>1</sup> Donostiako Informatika Fakultateko IXA taldekoak.

<sup>2</sup> Ikus horretarako: Sarasola K., Euskarazko softwarearen katalogoa. *Ehlyur. Zientzia eta Teknika*, 117. zk. 60-62 orr. 1997.

<sup>3</sup> Ikus horretarako: Artoia X., Informatika eta euskara: gaur egungo arazoak eta aurrera jotzeko bideak. *Utzaro*, 20. zk., 77-92 orr. 1997.

ordenadoreak urrun ikusten dituzten giza adimenaren ahalmen linguistiko orokorrak. Baina horrek ez du esan nahi hizkuntza lantzeko tresna automatikoak utopia direnrik, hizkuntzaren oinarriko ezagutza minimo batekin laguntza interesgarriak eskain daitezke eta. Testu guztiak ez dira zailtasun maila berekoak: ez da berdina ulertzea "Obabakeroak", telebistako eguzaldi-iragpena, edo egunkariko zinemakarteldegia. Hirurak euskarazko hizketa izan arren, bakoitzean erabiltzen diren hitzak, esateko erak eta esanahiak maila desberdinekoak dira erabat. Euskaldun alfabetatu batek ederto ulertuko luke hirur kasuetan, baina ordenadore bidez ulertu nahi duen programa batek zailtasun handiagoak izango ditu, lehenengo bi kasuetan behintzat. Emaitza probetxugarriak lortzeko, ordenadorearen lana aztergai espezifikoa eta mugatu batean kokatu behar da. Egun aurretiko hizordua ematen duen sistema gehienek zerbakiak eta astegunaren izenak besterik ez dituzte ulertzen, baina hala ere ekonomikoki oso interesgarriak diren aplikazioak antolatzen dira horrekin. Etorrizunean, aplikazio mugatuko sistemak bilduz, lor litezke ahalmen handiegoko sistema berriak, baina egun ibili dabilzan aplikazioek helburu espezifikokoak dituzte.

LNParen barruan azaltzen diren sistemak eta produktuek hobeto aurkeztearren komeni da bereiztea zein diren aplikagarritasun-maila desberdinak. Lau multzo nagusi egingo ditugu: lehenengoa, linguistikaz edo informatikaz gutxi dakten erabiltzaile arruntarentzat salgai diren **aplikazioak** sartuko ditugu; bigarrena LNPko ekoizleentzako bakarrrik interesgarriak diren **tresnak**, produktu berriak garatzeko baliagarriak; hirugarrena aztertuko ditugu edozein aplikazio edo tresna garatzeko behar-beharrezkoak izango diren **oinarririk**; eta, azkenik, laugarrena oraindik aplikazio mailara ailegatu ez diren **ikerketak-gaiak** sartuko ditugu.

Artikulu honetan LNParen egungo egoera bi ikuspuntutatik aztertu dugu. Hasieran merkatumen aurki daitezkeen aplikazioak aurkeztuko ditugu. Horietako gehienak ingelesaren munduan mugitzen dira eta bigarren maila batean frantsesa, alemaniera eta espainiera bezalako hizkuntzak dandue; euskararako aplikazio gutxi ditugu oraindik. Artikuluaren bigarren zatian helburu edo aplikazio horietara noizbait heltzeko beharrik gemituzkeen abiaburuak deskribatuko dira, epe erdi edo luzean euskararen tratamendu automatiko zabaldua posible egingo duen tresnak eta oinarriak, batik bat Donostiako Informatika Fakultateko IXA taldean azken hamar urte honetan sortu direnak<sup>4</sup>.

## Aplikazioak

LNParen 40 urteko historia gora-behera handiak izan dira. Helburu jilugarriak lortzeaz zeudela use zuten une euforikoaren ostean, belarririk jaisi eta helburu apal baina euskaragarrigarrietara mugatzeko une pragmatikoak etorri dira birritan edo. Erabateko itzulpen automatikoa ordenadorearen eskutik etortriko zela aurrekusi zuten 1954an Georgetown-eko Unibertsitate inguruan. Alabaina, 1966an itzulpen automatikorako diru-iturri ofizial guztiak itxi egin ziren, ALPAC txosten ezagunak horrela gomendatu eta gero. Aurrerago, 1980 inguruan, adimen artifizialeko teknika berrien eskutik ordenadoreak geure hizkuntzaz—lengoia naturalaren—programatuko gemituela agintzen zitzaigun.

<sup>4</sup> Artikulu osoan zehar aipatuko ditugunak talde horren esperientziari eta egungo jardunari loturakoak izango dira gehienbat. Informazio gehiago nahi duenak jo beza amarnumeko helbide honetara: <http://www.ji.sihnu.es/Groups/IXA>.

Gaur egun ahaztuta daude horrelako ametsak. Dena dela euforia eta pragmatismoko ziklo horiek bi motako emaitzak utzi dituzte: alde batetik, hoberio baloratu eta ezagutzen dugun hizkuntzaren egitura eta erabilera, eta aitortu behar izan dugun ez direla hasieran usie bezain sinpleak; bestetik, helburu utopiko horiek lortzeko asmoan eraldi diren tresnekin helburu apalagoa duen baina komertzialki bideragarriak diren produktu asko merkaturatu dira. Horrelako zenbait aplikazio arrakastatsu aipantuko ditugu ondoren.

### *Testuen edizioa eta gestioa*

Ordinadoreak kalkulu ugari eta konplexuak egiten dituen makina dela esan daiteke. Programak idazteko erabiliko zirela usie genuen, ez ordea prosa idazteko. Baina zer dela-eta zabaldu zait atea etxe eta bulego gehienetan? Erreko jaun-andreek programatzen ikasi dutelako? Ez, ezpada ordenadorea testuak idazteko idazmakina azkarra delako, edo Interneten bidez hainbat informazio eskuratzeko tresna ona delako. Edonork ulertu eta idazten du prosa baina guxiak programak. Testu-ediziorako eta testu-gestiorako tresnen garapena guztiz lotuta dago azken urteotan konputagailuen erabilera masibarekin. Honez gero testu-edizioa ez da tekleatze hutsa, edo testu baten bertsisio berri bat lortzea aurreko bat kopiatu, moldatu edo osatu, ezta hamaika formatu edo itxura desberdinetan aurkeztu ahal izatea bakarrik ere. Egun badira testu-egileari eskaintzen zaizkion laguntza bereziak. Ikus ditzagun orain zein diren garrantzitsuenak.

*Ortografia-zuzentzaileak* bete dituzte urte batzuk merkatuan, eta gaur egun hizkuntza askotarako aurki daitezke. Zuzentzaile hauek testuko hiz bakoiza aztertzen dute ea hiz posiblea den egiazatzeko, baina gehienetan testuingurua kontuan hartu gabe. Euskara bezalako hizkuntzen kasuan hitzak kasu desberdinetan deklinaturata agertzen direnez askoz lan konplexuagoa da hiz zuzenak eta okerrak bereiztea, analisi morfologikoa egin behar baita. Hala ere, 1994tik dago dendetan XUXEN euskararako egiazatatzaila/zuzentzaile ortografikoa. Zenbait hizkuntzarentzat *idatzkera- eta sinaxi-zuzentzaileak* ere merkaturatu dira; hauek testuingurua kontuan hartzen dute eta, adibidez, "nik joan naiz" esaldia prozesatuz gero, ortografia-zuzentzaileak ez luke erorerik salaruko, hiru hitzok isolaturata posibleak baitira, baina sinaxi-zuzentzaileak testuinguru horretan "nik" hitza gaizki dagoela salaruko luke eta "ni" izan beharke lukeela proposatu. Nahiz eta erore guztiak harrapatu ez, laguntza ederra eskaintzen diote eskutitzak edo bestelako txostenak idazten dituenari. *Laguntza lexikalean* edozein hitzen sinonimo edo antonimok lor daitezke testu-prozesaketako programatik atera gabe, baita taxonomikoki konkretuagoak edo orokorragoak diren antzeko hitzak ere (adibidez: *insektu* hitzetik orokorrago den *animalia* edo konkretuagoak diren *inurri, enli, ...), thesaurusa* konstitulatu.

Testu eleamntzak lantzeko adibide gisa Siemens-en Eurolang Optimizer, IBM-ren TranslationManager/2 eta Trados-en Translation Workbench programak aipatu behar dira. Prozesadore zabalduenetan integratzen diren programa hauek glosategi, hiztegi eta itzulpenen berrerabilenerako laguntzak eskaintzen dituzte. Itzulzaileek arazo franko izaten dute terminologiarekin testuko gaiarekin ohituta ez danduean. Horrelakoetan terminologi akatsak (hitz ezegokia ematea edo termino bera orritalde berean desberdin itzulitzea) askoz itsusagoak dira erore ortografikoak baino. Glosategi, hiztegi orokor edota hiztegi berezituoen on-line moduko erabilerak

erore horiek urritzeko baliabide eraginkorrak dira. Itzulpenen berrerabilpenerako laguntzek itzulzaileari lana errazten diote testuen bertsisio berriak egiterkoan, aldatu dena itzuli beharke bainu eta ez testu osoa. Automatikoki sortzen dute bertsisio itzuli berrira, osatu behar diren hutsuneak bereizita agertzen direla. Siemens-en EuroLang Optimizer programak Metal itzulpen-sistemarekin batera ere lan egin dezake.

Testu-masa handiak tratatzeko edo gestioatzeko aplikazio nagusiak lau dira: kontzeptu-bilaketa, kategorizazioa, informazio-erazketa eta testu-sofokuntza automatikoa.

On-line moduko *kontzeptu-bilaketaren* inguruan mila milioi dolarreko industria antolatuta zegoen 1994an Estatu Batuetan. Orain arte erabiltzeko teknika oso sinpleak ziren (hitz gakoan kombinazio boolear hutsa); gaur egun lematzazioa, perpausen butkaeren detektzioa, akronimoen zabaltzea eta kalkulatu estatistikoak ia sistema guztietan egiten dira; Clari, Conquest eta ConText (Oracle) produktuen egileek, etorkizuneko bidea erakutsiz, beren ekarpenaren iurrria LNPKo teknika sofistikatuenean kokatzen dute. Euskararako ere bada berriki Ametzagana taldeak kaleratutako Kapsula softwarea, euskarazko dokumentu-basen gestioa zuzendua.

*Kategorizazio-sistemak* oso baliagarriak dira makina bat dokumentu (adibidez: telefonoetako matxura-partreak, albisteak, hildako militarren partreak, marketineko datuak, ...) kategoriatan multzo txiki baten arabera sailkatu behar izanez gero. Esate baterako, Carnegie Group enpresaren Construe sistemak Reuter informazio-agentziaren artikulak automatikoki sailkatzen ditu, eta urez ure agentziari 750.000 dolarreko aurrezpena ekarri dio 1990 ureaz geroztik. ATT telefono-kompaniak daukan sistemak matxura-partreak automatikoki bideratzen ditu konponketaz arduratu beharke den bulegoraino. Zenbait aplikazioetan, nahiz eta dokumentu guzti-guztiak ezin sailkatu, nahikoa da dokumentuen gaineko estatistika orokorrak lortzea. Esate baterako, matxura-parteen estatistika orokorrak produktzio-kaie baten osagai ahulenen zein diren jakiteko balio dezake. Kategorizazio-sistema batzuk haratago doaz eta satatzen dira ezagutzen zein diren "elementu agerue berriak", hau da, nahiz eta sarritan azaldu, behin eta berriz sailkatu gabe geratzen direnak berantiarzako kategoriarik definitu ez delako. Dokumentu-kategorizazio teknika bi motakoak dira: estatistikoak eta ezagutza-ingeniartzakoak. Ezagutzaren ingeniartzako teknika baliatzen dituztenak zehatzagoak dira, kalitate hobea lortzen dute, baina oso garestiak dira, eta ez dira errentagarriak dokumentu kopurua benetan erraldoia ez bada (horrela aitortzen zuten Carnegie Group enpresakoek). Dena dela, tresna estatistikoak erabili behar badira ere, aldeaz aurreik testu luzeak etiketatu behar dira sistemak horietatik ikas dezan.

*Informazio-erazketerako sistemak* lengoia naturaler idatziriko testuetatik datu-base egituratu bat osatzen dute. Azken helburua albiste multzo handi batetik abiatuz fitxa konkretuak bereizte litzateke non-nori-zer egin dion jakiteko. Dena dela, helburu apalagoak baldin badituzte ere, badira produktu asko merkatuan euekin handiak ateratzen dituztenak. Gehienek dokumentuak aztertzen dituzte enpresa, pertsona, hitzordu-data, telefono edo zerbizuen erreferentzia hutsen bila. Adibidez Westlaw eta Lexis-Nexis-ek enpresen aipamenak bilatzeko programak saltzen dituzte, enpresaren aipamena modu askotara azal daitekeelarik: esate baterako, IBM, I.B.M., International Business Machine, e.a.

*Testu-sorkuntza automatikoa* informazio-erazketaren kontrako da. Kasu honetan ordenadore barruan dauden datu konplexueatik abiatuz (inprimakiak, datu kodeatuak edo zenbakitzko formatuan dauden informazioak...), datu horien edukia azalduko zaio erabiltzaileari bere hizkuntzan. PLANDOC sistemak telefono-enpresa batenzat honen telefono-sarerako hobeakuntzak asmatzen ditu, baina erabiltzaileak ez du ikusien programaren emaitzaren kode uterrezina, berori azaltzen duen ingelesezko testua baizik. Forecast Generator sistemak (geroago aipatuko den METEO sistema ospetsuaren ondorengotzat hartua izan dena) ingeles edo frantsesezko testuak idazten ditu ordenadore batek kalkulazten dituen eguraldi-iragapen kodeetatik abiatuz.

### *Hizulpen Automatikoa*

Produktu ugari dago merkatuan salgai testu-izulpenean laguntza emateko, baina euskara tratatzen duen sistemarik ez dago. Hizulpen perfektua egiten duen sistemarik ez dago inon, eta sistema batek berak ere ez ditu testu hierarioak izultzen. Guztirik izulpen teknikoak dute erabileremu, testu teknikoetan hizkuntzen arteko hitzen eta esaldien korrespondentzia anbiguotasun gutxiago aurkitzen baita. Sistemaren batek tratatzen dituen testuak edozein motakok badira, ziur emaitzaren kalitatea kaskarra dela. Hala ere, gero azalduko dugun bezala, zenbait kasutan sistema horiek lagungarri izango dira.

Izulpenaren automatizazioa ez da ia inoiz erabatekoa, eta automatizazio-mailaren arabera ondoko sailkapena egiten da: 1) Erabateko izulpen automatikoa: errealtatea baino ametsa da gaur egun, non eta helburua ez den edukia iren ideia orokorra ateratzea. 2) Giza laguntzaz burutuiko ordenadore bidezko izulpena: lanaren gidaria makina da, baina fase desberdinetan laguntzak eska ditzake; hitz baten adiera zuzena hautatzeko edo esaldi baten analisi nondik hasi behar den erakusteko adibidez. 3) Ordenadorez lagunduriko giza izulpena: lanaren gidaria pertsona da, baina ordenadoreaz baliatzen da hiztegi berezitan kontsultak egiteko, testuaren formatua txukuntzeko eta zailtasunik gabeko testu-zatiak izultzeko. Agian izulpenaren zati handi bat ordenadoreak egingo du ia laguntzarik gabe, baina testua egokitzeko aurreprozesaketa edota emaitza zuzentzeko postedizioa ohikoak izaten dira. 4) Datu-banku terminologikoak: hiztegi berezituak erabiltzeko aukera hutsa eskaintzen duten laguntza-sistemak.

Testu izulpen erabileraren nagusi bi bereizten dira: edukia iren ideia orokorra ateratzen dutenak, eta zabalkuntza handiko informazio zehatzak izultzen dituztenak. Lehenengoren adibide tipikoa "internautarena" dugu: hizkuntza arrotz batean ixurra oneko web-orri bat aurkitu du, guztia zehatz-mehatz izultzea denboraz edo diruz oso garestia litzaieke, eta gainera, ia ziur oso-osorik ez litzaiokeela interesatuko gero. Guztiz zuzena ez den baina merkea den izulpena gantbegiratu jakin ahal izango du benean interesatzen zaien partea zein den, eta gero zati horren izulpen zehatza lortu. Oro har, denbora eta dirua irabaztiko du honela. Beste aldetik, zabalkuntza handiko informazio zehatzen adibide gisa, etxetresna elektronikoko baten erabilererako azalpenak ditugu. Testu horien zehaztasun-uleragarritasunak produktuen arrakastarako gilitza izango dira. Beraz, izulpenak kalitate handikoa izan beharko duenez nahitaezko lana izango da giza-izultzaile baten zuzenketa edota postedizioa. Postedizioa ekiditeko asmoz zenbait sistematan satatu dira jatorrizko testuak mugatzen, erraz itzuli ahal

izango denera mugatuz. Horrelakoetan analizatzaile bereziak definitzen dira jatorrizko testuetan lengoia kontrolatutik ateratzen diren hitzak edo esaldiak salatzeke.

Montrealeko TALUM taldeak egindako METEO sistema da emaitzarik arrakastatsuenen lortu dena. Pate meteorologikoak izultzen ditu 1977tik hona, ingelesezik frantsesera, eta izulpenaren %80 erabat zuzena da. Eguneroko oso antzekoak ziren izulpen aspergarri hauek egiteko izultzaileak bilatzea zaila zen, nahiz eta soldata ederrak eskaini. Ute hartatik hona lana eguneroko burutzen da METEOren laguntzaz. Hamaita saio egin da geroztik sistema honen diseinua beste gai batzuetara zabalteko, baina ezin izan da horren birbilua den beste gai bat aurkitu. TALUM taldeak berak hegazkinetarako eskuliburuak izultzeko saiok egin zituen, baina hasierako emaitza itxaropenisuek piztutako ametsak laster itzali ziren.

SYSTRAN Institutuua 1970. urteaz geroztik izulpen automatikorako tresnen saltzaille nagusia izan da. NASA, Europako Batasuna, General Motors eta Xerox dira bere bezorik ezagunenak. Europako Ekonomi Elkartetik egokitzapen nekeatsua behar izan zuten —100.000 hitzeko hiztegia definitu behar izan bait zuten— frantses/ingeles izulpena ahalbidetarako. Baliagarritasun-mailaren berr emateko edo, aski izan daitieke aipatzea nola orain dela hamar urte 20 izultzailek erabiltzen zuten sistema hau Luxemburg-en, hilabetean milaren bat orrialde ingeles/frantses, frantses/ingeles eta ingeles/italiera bikoteetarako izultzen zutelarik. Kanadako General Motors-ek eskuliburuak izultzen zituen ingelesezik frantsesera: 130.000 hitzeko hiztegia definitu ondoren, izultzaileen lana lehen baino 3 edo 4 aldiz arinagoa zen, eguneko 1000 hitzetara helduz. SYSTRANen oinarri informatikoa, baina, guztiz atzeratuta dago. 1960ko hamarkadako teknologia erabiltzen baitu; hala ere, Systran izulpen-sistema hoberearen eta erabilienaren artean dago oraindik.

Dozenaka produktu dago ordenadore pertsonaletan izulpenak egiteko hizkuntza bikote desberdinen artean. Adibidez, ingeles-espaniera izulpenak egiteko badira sistema komentzialak: Spanish Assistant, Dos amigos, Context, Translate, Globalink. Guztietan postedizioa beharrezkoa da, eta nolabaiteko elkarrekintza dago beti giza-izultzailea eta programaren artean hitzen adiera zuzena hautatzerakoan eta. Bizkaiko GENISA empresa ere ingeles-espaniera bikoterako sistema bat garatzen ari da, eta euskara ere lanzen dute maila apalago batean.

### *Ordenadoreen erabileraren Lanaren bidez*

Aplikazio-mota honetako sistemak, ordenadore eta gizakiaren arteko komunikazioa errazten dute, erabiltzaileak bere hizkuntzaz lan egiteko aukera du eta. Horrelako sistemak inplementatzen zaitiak dira; galdera eta erantzuneez osatutako elkarriketa ulertu ahal izateko, partaideen planak eta helburuak aztertzeke tresnak behar baitira. Hiztun bakoitzak momentu bakoitzean zer dakien eta zer nahi duen asmatu behar da eta, gainera, ezagumendu horiek elengabe eguneratzen ibili behar da elkarriketa aurrera joan ahala. Helburu orokorretik ez da luzarotan salgai egongo, baina badira dagoeneko aplikazio konkretuetan lortu dauden batzuk.

Datu-baseetarako galdeketa-sistema ugari dago, batez ere ingelesez. Datu-base konplexuei galderak egin ahal izateko lengoia berezi bat ezagutu beharrik datu-baseen erabiltzaile potentzialen

kopurua murrizten duenez, galderek lengoia naturaler egin ahal izatea oso interesgarria da bezero berrak harrapatzeko. Horrela produktura lehenengo aldiz hurbiltzen den erabiltzaile potentzial berrak oztopo guxtiago aurkituko du martxan ikusiteko. Behin produkturaren funtzionamendua ezagututa motibazioa etorriko zaito probetxu handiago ateratzeko eta, horiaz, kontsulta-lengoia berezia ikasiteko. Izan ere, ordurako amna janda dauka eta programaren munduan sartuta dago.

Symantec-en "Question & Answer (Q&A)" sistemak arrakasta ederra izan du 1986 urtez gero. Sistema hau analisi sintaktikoa bigarren mailarako lagatzen duten "gramatika semantikoeetan" oinarritzen da. Galderek oso zailak ez badira, emaitza harrigarriak lortzen ditu.

Alde ikaragarria dago ordenadore erraldioetarako eta mikroetarako egindako interfazeen artean; bai prezioz (milioiak eta hamamaka mila pezeia inguru, hurrenez hurren) eta bai ahalmenez. Hizkuntzaren tratamendua askoz zabalago eta sakonagoa izateaz gain, sistema handietan erabiltzailearentzat laguntza eta erraztasuna handiagoa da. Erabiltzaile antizi erantzun dakitoke eta datu-base ahalsuagoak arizitzeko aukera eskaintzen dute. Mikrooan kokatuako interfazeak guztiz desberdinak dira. Merkatuan 100 baino gehiago dira ingelesezko produktuak.

Ikerkuntzaren munduik datozen ahalegin berrietan LNPKo teknirik eta multimedialakoak bilzteko saioak egiten dira, edo menuen bidez erakusten dira egin daitezkeen esaldien egitura eta kontzeptuak. Adibidez Texas Instruments enpresaren "Natural Link" paketean, erabiltzaileak ezin du galdetu edonola, berari aurkezten zaiou menu moduko pantaila batean hitzak edo esaldi-zatiak hautarzen ditu nahi duen galdera osatzeko. Horrela esanda, menu hutsa dela dirudi, baina bere azetik dagoen hizkuntz analizatzailea antzeko beste sistemem mailakoa da. Pakete honen ezaugarririk onena gantentasuna da: erabiltzaileak ondo dakit zein izan esaldi ulertuko diren eta zein izan ez.

Datu-baseei buruz azaldu dugun hori guztia bertin esan daiteke gainontzeko aplikazio-programez ere. Gehienak adimen artifizialeko sistemetan integratuta daude. Baina bestelakoetan ere badira eta, adibidez, zenbait konzeptu-bilazaitetan galderek lengoia naturaler (ingelesez kasu guztietan) egin daitezke.

### **Ahozko hizkuntzaren tratamendua**

Ahozko hitzak edo esaldiak ulertzea zaila da, hizkuntza idatzia ulertzeko arazoei ahozko hizkuntzaren problematika erantsien zaitolako: hitzak ez dira guztiz bereziz egin egiterakoan, esaldien hasiera eta bukaera erdikoa baino inentistate txikiagoz ematen dira, eta, gainera, seinale fisikoan zaratak ohiko oztopoak izaten dira.

Sistema gehienek oso hitz gutxi ezagutzen dute, eta horien artean beti daude zenbakiak. Horrela erabiltzaileak zenbaki bat ahoskatuz aukera desberdinen artean hautatuko du behin eta berriz menu desberdinetan zer (edo zer eskatu) nahi duen ondo zehaztu arte. Merkatu handia zabaldu da horrelako sistemak telefono bidezko zerbizuetan integratzeko: aurretiko hizordua, produktu-eskerrak, e.a. Beste alde batek, hizketaren ezagutzarik gabe, gero eta arnuntago bihurtzen ari zaigu makinem ahots sintetizatuek entzutea gasolindegietan edo tabako-edariak saltzen dituzien makinetan.

Natural Vox enpresa arabarrak aurretiko hizordua —medikuraren eta errenta-aitorpena egiterakoan— automatikoki lortzeko sistema telefonikoak ezarri ditu azken urteetan, eta arrakasta handia izan du.

Ahozko hizkuntzaren tratamenduko teknikak antzeko beste aplikazioetan ere erabiltzen dira: eskuz idatzitako testuak ezagutzeko edota testu mekanografiatuen bertsiio elektronikoa lortzen duten OCRetan (Optical Character Recognizer, karaktere-ezagutzailer optikoak).

## **Abiaturrak**

Batez ere ingeleserako merkatuan aurki daitezkeen aplikazioak ikusi ondoren, artikularen bigarren zatian helburu horretara noizbait helduko bagara martxan jarri behariko genituzkeen abiaturrak deskribatuko ditugu, beti ere, IXA taldean markaturako estrategiarri jarraituz eta bere mugekin, noski, eta azken hamar urte honetan sortu ditugun tresna eta oinarrietatik abiatuz. Ahiaburuen artean aipatzekoa lizateke, jakina, artoko ikertuntza. Hala ere, artikulu honetan aplikazio eta tresnen oinarri direnak azalduko ditugu bank bat, eta egun lanitzen ari garen bestelako ikerketagaitz —teorikogoko edo— aipatu baino ez ditugu esingoo azkeneko aralcan.

### **Tresnak**

Ahal honetan hizkuntzaren tratamendurako aplikazio-ekoizleentzat edo artoko ikertzaileentzat interesgarriak diren tresna batzuk ikusiko ditugu. Tresna horiek ez daude diseinatutrik, oro har, erabiltzaile arruntarentzat.

### **Analizatzaile morfologikoa**

Ingelesaren flexio-morfologia simplearen eraginaz ordenadorez egindako analisi/sintesi morfologikoiari kasu handiegia ez zitzaion egiten, eta askotan aplikazioak ibiltzen ziren hizkuntzaren forma guztiak zintuen hizitegi batekin. Hau pentsaezina da euskara bezalako hizkuntza flexionatu eta eranskariaren kasuan, erro batek sor daitezkeen hitz flexionatu posibleak asko eta asko baitira. Eta hori ez bakarrik euskararako, beste hizkuntza askotan (suomiera, turkiera, etab.) arazo bera baitzegeen. LNPKo teknikak ingelesetik beste hizkuntzetara hedatzean eman zitzaion morfologiari daukan garrantzia.

Morfologia automatizatzeke orduan hizkuntzaren hiru aspektu deskribatu behar dira zehatz: lexikoa, hizaren osaketa eta aldaketa fonologikoak. Lexikoa funtsezko elementua da eta lexiko oso eta orekatu bat eraikitzea lan izugarria da aurretik datu-base lexikal bat ez badago, lexikoko sarrera konbentzionaler gain morfema ez-independenteak ere behar dira eta gainera osagai guztiei dagokien informazio morfologikoa. Hizaren osaketa (morfotaktika edo hizaren gramatika deitu ohi da, hizkuntzak onartzen duen flexio-bideari jarraiki sarrera bakoitzaren ondoren etor daitezkeen elementuen deskribapen formalak) deskribatu behar da sistema informatikoak jakin dezan nola lot daitezkeen erroak, aurritziak eta atzizkiak. Horrela lortuko dugu ez ezagutzera —eta ez sortzea—

*erxego* baina bai *handiago*, eta ez *batilako* baina bai *batia* eta *delako*. Azkenik, elementuak biltzean sortzen diren aldatuketak (aldaketa fonologikoak) deskribatu behar dira; horrela azalduko da adibidez, aurretik azalduetako adibideari jarraituz, *bat* eta *da* biltzean ez dela *batida* gertatzen, *batia* baizik.

Hiru atal horien bidez egien dira gaur egun hizkuntzen deskribapenak morfologia automatikoari begira. Kasu gehienetan flexio-morfologia hartzen da kontuan baina ez eratorpena eta elkarketa, azken horiek ez baitira erregularak.

Deskribapen horiek egiteko modua erabiliko den programaren menpe egongo da (programak hizkuntzari berezi egien dira gaur egun). Horrekin lor daiteke analizatzaile/sortzaile morfologiko bat, programa hauek askotan gai izaten baitira hitza emanda analisisa lortzeko eta erro baretik abiatuta deklinabide osoa lortzeko. Horixe izan da IXA taldean euskararako egin dugun lehen tresna.

Analizatzaile morfologikoa oinarri bat da hainbat aplikazioetarako. Honako hauek dira garrantzitsuenak:

- Zuzentzaile ortografikoa. Hitza analizatzerik baldin badago hitza zuzena izango da eta bestela abisu bat emango da. Hori egingdu dago euskararako.
- Tutore-sistema automatikoak hizkuntza ikasten ari den jendearentzat. Erroreak detektatzeko, arketak programatzeko etab-erako oso elementu interesgarria da. Honean ari gara lanean gaur egun.
- OCR dokumentuen irakurketan (eskanerrak erabiltzean) sor daitezkeen erroreak detektatzeko.
- Hizketaren sintesia edo testu-sorkuntza lortzeko sorkuntza morfologikoa funtsezko osagarria da.
- Hizkuntz aplikazio sofitikatuagoetarako —sintaxian oinarrituakoa, itzulpen automatikoa, etab.— lehen urrats gisa.

### ***Lematizatzaile/etiketatzailea***

Morfologiatik sintaxira doan bidea oso luzea gertatzen da LNParren munduan, morfologia-sistema osoak erakitzera posiblea den bitartean, gaur egun ez baita oraindik posible sistema sintaktiko automatiko oso bat garatzea. Are gutxiago ingelesa bezala ikertu ez den hizkuntza baterako. Hori dela-eta tarterko bideak hartu dira eta analizatzaile sintaktiko orokorrak baino sinpleagoak diren tresnak bultzatu dira azken urteetan. Arrakastatsuenak etiketatzaileak izan dira eta, hakekin batera, lematizatzaileak. Etiketatzaileek testuko hitz bakoiak dituen analisi desberdinen artean zuzena dena aukeratu behar du; lematizatzaileek, aldir, lema posibleen artean dagoakiona. Adib. *zuen* hitza analizatzean posible da *ukan* aditza lehen aldiari izatea edo *zuek* izenordaina genitiboan. Testuinguraren arabera erabaki behar du lematizatzaileak zein den hitzari dagoekin etiketa zuzena (aditza edo izenordaina) edota zein den bere lema (*ukan* edo *zuek*). Beraz, lan hau konplexua da, ez baita posible hitz isolatutak aztertzea, eta sintaxiaren lehen urrats gisa hartzen da tresna hauena lana.

Aral nagusia desambiguazioa bada ere, beste zereginak ere badaude halako tresna bat garatzean, esate baterako, hitz antzezko unitate lexikalen identifikazioa (lokuzioak, hitz-elkarkeak, pertsona-izen osoak, etab.). Desambiguatzeko teknika bi ildo nagusitanik doaz: metodo empirikoak edo estatistikoak baretik, eta ezagumendu linguistikoan, erregelatan, oinarritutako metodoak bestetik. Gero eta joera gehiago dago bi metodo-motak konbinatzeko. Lortzen diren emaitzak ez dira zeharo fidagarriak baina %95-98 tarteko fidagarritasuna lortzen da. Sistema hauek garatzen ari dira hizkuntza askotarako eta gu euskararenta ere egiten ari gara.

Tresna hauek izan duten arrakasta beren aplikazioetan datza, oso aplikazio interesgarri eta aktualak baitintuz lematizatzaile/etiketatzaileek:

- indexazioa: testuak indexatu nahi direnean ez zaigu forma interesatzen, lema eta kategoria baizik. Indexazioa da oinarria gaur egun hain modan dauden datu-base dokumentaletan eta Interneteko bilatzaileetan. Adibidez, testu batean *kaltekoak*, *kalera* eta *kalajirarik* agertzen badira, lehen biek azaldu behar duie *kalera*z galdetzen dugunean, baina hirugarrenak *kalajiraz* egien dugunean.
- terminologia/lexikografia: automatikoki lemak ondo identifikatzen badira eta dagozkien etiketak egokitzten bazaizkie lan lexicografikoa erruz errazten da, eta testu baretik terminologia automatikoki erazteak ez dirudi oso lan zaila.

### ***Analizatzaile sintaktikoa***

Analizatzaile sintaktikoen zeregina testuetako osagai sintaktikoak ezagutzera da: hitz isolatuez osatu sekuentzietan elkarrekin lotuta dauden egitura sintaktikoak (perpausak, izen-sintagmak, aditz-sintagmak, izen-lagunak, eta abar) ezagutuko dira. Analisisaren oinarria lexikoa eta gramatika izango dira, hizkuntzako hitzen ezaugarri sintaktikoak eta egitura sintaktikoen osaketa posibleak definituko dituztenak.

Prozesu honek ambiguetate handia sortzen du, esaldi bakar baterako analisi posible anitz lor baitaitezke. Kontuan hartu, gainera, hitzen analisi morfologikoa aldez aurretik egin behar dela, eta hitz bakoiatzeko analisi morfologiko anitz sortzen direla (euskararen kasuan eta gure datuen arabera, hitzeko 2,7 aukera desberdin batez beste). Ambiguetatea eragozpen handia da erabateko analisi automatikoa lortzeko. Erraza da laborategiko esaldi multzo bat prozesatuko duen analizadoreara erakitzea, baina oso zaila testu libreekin lan egingo duena egitea. Ingeleserako ALVEY sistemaren gramatikak edozein esaldi tritatzen duela diole, baina gero ez da oso erabilgarria beste analisi edo sistemetan aplikatzeko; darabigun hizkuntzaren ambiguitasuna dela-eta testu arruntetako esaldietan batez beste 100 bat analisi desberdin lortzen baitute.

Formalismo asko dago gramatikak definitzeko, baina gehienak Chomsky-k definitutako Testuinguririk Gabeko Gramatiken gainean egindako hedapenak dira. Baterakuntza-gramatiketan erregela bakoiatzeko osagaien gainean ekuzio-multzo bat definitzen da osagaien gaineko komunizadura egiaztatzeko eta egitura komposatutak osatzeko. Beste planteamendu berri bat agertu da zenbait sistema bertian: xedea ez da esaldi oso-osoa analizatzea, hori bilatuz gero gehienetan porrota izango baita

emaitza, eta nahikoa da esaldiaren azaleko analisia lortzea, hau da, bereiztea zein diren osagai posibleak eta beren arteko loturak. Emaitza hauek teoria linguistikoan ikuspenetik ez dirudite oso dotore, baina emaitza horiekin beste hainbat tratamendu informatikorik atea zabaltzen zaio. Beste alde batek, azken analitzaile horiek konputazionalki askoz azkarragoak dira. Murrizpen Gramatikak dira adibide ezagunena.

## Oinarriak

### *Datu-base lexikala eta morfologiaren deskribapena*

Datu-base lexikala da hizkuntzaren lexikoaren biltegi erraldoia. Hiztegi elektronikoko moduko bat da, hizkuntzaren tratamendu automatikoa begira eraikia, eta, beraz, hizkuntzaren tratamendua automatizatu nahi horrek diuen eskakizunak kontuan harturik antolatua. Horrek eskatzen du, noski, lexikoaren antolakuntza gero zertarako erabiliko den kontuan hartuz egitea, eta lexiko-deskribapenaren sistemazioa bat: sarrearen kategoria-sistema bateratu eta homogeneoa erabilizera, kategoria bakoitzeko elementuak behar den bezala deskribatzeko beharrezko diren ezangariak zehaztea, etab.

Euskararen kasuan, IXA taldean Xuxen ortografia-zuzentzailearen prestatze-lanari ekin genionean sortu zitzaigun halako lexiko-biltegiaren premia. Gorago esan bezala, baina, zuzentzaile hori oinarritzailekoa zen analitzaile morfologikoaren azpiprodukturatz hartzen genuen guk, eta datu-base lexikala ere ez genuen antolatari nahi izan zuzentzaile horretarako hiztegi edo hitz zerrenda soil gisa, etorkizunean euskararen tratamendu automatikoaren arloko beste edozein tresna edo aplikazioarako oinarri lexikal sendo gisa baizik. Eta horrela sortu zen EDBL, Euskararen Datu-Base Lexikala, harez gero gure lanerarako oinarri lexikala izan dena, etengabe eguneratuz joan dena, eta gaur edo bihar komunitate zabalgogo bati bere atek irekiko dizkiona, oinarriak prestatze-bide honetatik beste batzuk ere baliatu daitezzen.

Datu-basa diseinatzerakoan garrantzi handia eman zitzaion, bada, etorkizunean izan ditzakeen hedapenak onartzeko behar bezain malgua izateari eta, bereziki, bertan jasoko zen informazio linguistikoa ahalik eta erarik neutralenean deskribatzeari, hau da, formalismo edo teoria linguistikoerarik ahalik eta erarik independenteenean.

EDBLK gaur egun 70.000 sarreara inguru biltzen ditu, hiru atal nagusitan sailkatuta: hiztegi sarrearak (izenak, adjektiboak, aditzak, e.a.), adizkiak (aditz-forma jokatutak) eta morfema ez-independenteak (atizki, aurriki, e.a.). Sarreara-kategoria bakoitzeko aldeez aurretik definiturik dauden ezangari edo atributuak erregistratzen dira, eta kasu guztietan, lehen aipatu bezala, sarreari dagokion morfologia ere deskribatzen da (informazio morfotaktikoa), horretarako bi mailarako formalismoaz baliatuz<sup>5</sup>.

<sup>5</sup> Koskemiemi K.: *Two-Level Morphology: A general Computational Model for Word-Form Recognition and Production*. Ph.D. thesis, Helsinki University. Pubs. no. 11. 1983.

EDBL egun datu-basen kudeaketarako sistema baten penan dago eta halako sistemek ohikoak dituzten erraztasunak eskaintzen dizkio hizkuntzalariari —hizkuntzalariak baitira bere erabiltzaile nagusiak—: interfaze aisekina lanerako, informazioa egunean mantendu eta berorren konsistentzia bermatzeko erraztasunak, behar den aplikazioerako informazioa behar bezala iragazteko aukerak, eta abar. Euskararen baturatze-bidean izandako azken getakariak —Euskaltzaindiaren erabakiak, batik bat— egunean mantentzeko ere ezinbesteko tresna bihurtu da datu-basa, eta etorkizunean nahi genukeen zabalkundera izan dezarean EDBLK bete dezakeen lan importanteariko bat izan daiteke azken erabakien berri emango duen tresna izatea.

### *Hiztegi elektronikoa*

Hizkuntzaren datu-base lexikal orokorra oinarri dela, horren inguruan biltzen ahal dira beste zenbait tresna lexikal ere: definitzio-hiztegiak, hiztegi terminologiko berezitutak, hiztegi elebidunak, eta beste. Horrelakoen garrantzia ere ukatu ezina da, batez ere hizkuntzaren semantika tratagai denean edota itzulpenaren arloko aplikazioak egiterakoan.

Gaur egun, eta hementik aurrera zer esanik ez, ia ateratzen diren hiztegi guztiak ateratzen dira euskarri elektronikoren batean (CD-ROMean, batez ere), eta horietaz baliatzea ere helburu da hizkuntza baten tratamendu automatikorako oinarri lexikala prestatzeko orduan.

Gurean hor ditugu UZEIren EuskalTern datu-banku terminologikoa, I. Sarasolaren Euskal Hiztegia, eta Elhuyarrek, Harluxet Fundazioak eta Adorez taldeak, besteak beste, kaleratutako hiztegi-lanak euskarri elektronikoko desberdinetan; etorkizunean, eta datu-base lexikal zentral batekin behar bezala loturik, adibidez, hainbat lanerarako oinarri lexikal osagarri bihur litezke lan horiek, nahiz eta hasieran helburu horrekin sortu ez.

IXA taldean, oraintxe, Euskal Hiztegiarekin eta Aulestiaren ingelesa-euskara makinatutako bertsioarekin ari gara lanean, EDBLri lotuz honek oraindik ez duen osagai semantikoa (definitzioak) eta itzulpenezkoa (ingelesa) gehitzeko asmoan.

### *Gramatika konputazionalak: sintaxiaren deskribapena*

Sintaxia ere funtsezkoa dugu hizkuntzaren tratamenduren arloko edozein lani ekitiko, hizkuntza ezagutzea nahiz sortzea delarik helburua. Hizkuntzaren gramatika formalizatu eta konputazionalki tratatzeko moduan adierazi behar da, morfologiaz harantzago joan nahi duen edozein aplikazio edo tresnatan erabiliko bada.

Euskararen kasuan morfologia eta sintaxiaren arteko lotura estua hartu behar da kontuan lehenik. Horrek eraman gaitu tratamendu morfosintaktikoa analitzaile morfologikoa —morfosintaktikoa, hobeto esanda— integratzera. Hirzaren barruko gramatika modu bat definitu da horrela, eta, horri esker, analisi morfologikotik beretik prestatzen zaio bidea gero etorriko den analisi sintaktikoa. Ertegeela multzo baten bidez deskribaturik daude, bada, hitz barruko morfemen arteko erlazioak, erlazio horietatik hitz osarean analisierentzat garrantzizkoa den informaziorik eratortzen denean.

Baina horretaz aparte, euskarazko perpausen joskera, sintaxia, ere deskribatu beharra dago. Lehen esan den bezala, analisi zein sorkuntza sintaktikok ezinbesteko tresnak baititugu aplikazio gehienean. Horrela bada, beste hizkuntza batzuetarako baliagarri suertatu diren formalismoak eta analisi-teknikak erabiltzen ari gara gu ere. Horien artean azpimarratu nahi genuke Murritzpen Gramatika<sup>6</sup>, analisi morfologetikoa ateratzen den anbiguotasun maila jaisteko eta esaldien analisi sintaktiko azalekoa egiteko erabiltzen ari gara. Edozein testu—mugarik gabe—analizatzea helburu duen formalismo horretan hitzen "hurtileko gramatika" deskribatzen duten 1.000 erregela inguru idazriak ditugu oraingoz. Horretaz gain, PATR-II izeneko baterakuntza-formalismoz<sup>7</sup> euskarazko izen-sinagema eta perpaus bakunen egitura deskribatzen duen gramatika konputazionala ere garatu dugu.

### *Taxonomia semantikoa*

Hizkuntza ulertzea xede denean, baina, ez da aski morfoloia eta sintaxiarekin, semantikaz ere jakin behar izaten baitu programak. Anbiguotasun linguistikoa ebazteko beste modurik ez dago, asko eta askotan, semantikaz baliatzea baino.

Hizkuntza baten tratamendurako azpiegituran osagai semantikoa ere behar du bere lekua, beraz. Semantika lexikala da, beharbada, osagai horren prestakuntzan landu beharreko estreinatiko alderdia. Semantika lexikala hitzen semantika biltzen du, lexikoko elementuen artean dauden erlazio lexiko-semantikoa: sinonimia, antonimia, hiperonimia/hiponimia (klase/azpiklase erlazioak), erlazio meronimikoa (zaita/osoa, osagaia/osoa, e.a.), eta beste. Geroago etorriko da esaldien adierazpide semantikora eramango gaituen analizatzean.

Hiztegi aruntetan hitzen semantika lexikalari buruzko hainbat eta hainbat informazio dago. Informazio hori bertatik erazuteko—erdi-automatikoki, askotan—makina bat saio egin dira, baita gurean ere. Lan horien helburua, gehienetan, lexikoko unitateen artean erlazio lexiko-semantiko horiek esplizituki errepresentatzea izan ohi da, azkenik sare semantiko moduko bat lortzeko. Ingelesezko sare semantikoen artean ezagunena-edo WordNet izenekoa genuke<sup>8</sup>, eta euskararako halako sare bat eraztea genuke guk ere geure epe ertaineko jomngen artean.

Hitzen forma ezagutzetik harantzago joango litzatekeen testu-zuzentzean batek ezinbesteko luke halako tresna bat, eta gauza bera euskaraz idatzitako testuetan informazio-biaketa egingo lukeen hizkuntzari zuzendutako tresna batek. Adiera mailako desanbiguzioan ere beharrezkoa litzateke halako sare semantikoa.

6 Karlsson F., Vuolteen A., Heikkilä J., Anttila A. *Constraint Grammar: Language-independent System for Parsing Unrestricted Text*. Mouton de Gruyter, 1995.

7 Shieber M. *An Introduction to Unification-Based Approaches to Grammar*. CSLI Lecture Notes, no. 4, 1987

8 Miller G. *Five papers on WordNet. Special Issue of International Journal of Lexicography* 3(4), 1990.

### *Testu-corpusak*

Ikerkuntza-ardo honen azpiegituran nahitaezkoa den beste elementu bat testu-corpusak ditugu. Testu-corpusak testu-masa handiak dira, informazio linguistikorekin iturri nagusietako bat eta gorago aipaturako aplikazio, tresna eta oinarrietarako probaleku ezinbestekoak. Hizkuntz corpusak lexikografian duen garrantzia ezaguna den bezala, LNPrako lexikoi bat—datu-base lexikal orokorra bera—prestatzerikoa ere premiazkoa dugu. Gramatika bat ezin da hutsetik asmatu: testuak ditugu hizkuntzaren erabileraren lekuko. Egingdako tresnak eta aplikazioak ezin dira probatu laboratoriotiko hitz, perpaus eta esaldiekin soilik: testu errealak behar dira, egingdako programa horiek gero, benetako testuei aurre egin diezaietenean, porrot egingo ez badute.

UZEL eta Euskaltzaindaren EEBS corpusa<sup>9</sup> genuke, guk dakigula behintzat, egun Euskal Herrian lengoia naturalaren prozesamenduko lanerako dagoen corpus sistematzatu bakarra. Baina corpus hori, 3.000.000 hitzekoa izanik ere, ez da nahikoa. Testu-corpus horien bilze-lan eta antolatzea sistematioki ekin egin behar zaio lehenbaiten, modu planifikatu batean. Lan horretan toki askorako jendeak hartu behar luke parte—Euskaltzaindia, UZEL, komunikabideak, argitalerxeak, eta abar—uste baititugu halako lan bat behar-beharrezkoa dela, honetan ari garenontzat ez ezik, baita beste ikertzaile askorentzat ere. Gaur egun ez da duela urte guxi batzuk bezala, testuak euskarri elektronikoen egunero sortzen baitira, piñaka. Kontua da horiek biltzea, txukuntzea, eta ikertzaileon eskura jartzea.

### **Ikerkuntza**

Bukatzeko, lengoia naturalaren prozesamenduren arloko ikerkuntza hartuko dugu hizpide, izan ere ikerkuntza baita hizkuntzaren tratamendu automatikoa helburu duen programa ororen oinarri ezinbestekoa. Baina, gorago esan bezala, oraingoan aipatu baino ez ditugu egingo gure taldean gaur egun esku artean ditugun proiektu eta ikergaiak.

Morfologia konputazionalaren alorrean, murritzpen gramatiken bidezko desanbiguzioa eta estatistikari oinarritutakoa ditugu aztergai. Sintaxiaretan, berriz, sintaxian oinarritutako zuzenketa, gaur edo bihar gramatika-zuzentzaile bat egitera eramango gaituena. Aipatu bi arloak ukitzen dituela, Komputagailuz Lagunduriko Hizkuntzen Irakaskuntzan ere aurreratu da tes-lan bat.

Euskal aditzen azpikategorizazioa da beste ikergai teoriko bat, sintaxia ez ezik semantikaren artoa ere ukitzen duena, eta etorkizunean aplikazio praktikoa handikoa izatea espero duguna. Semantika eta lexikoaren alorrean, berriz, semantikan oinarritutako zuzenketai dituen arazoak azterzea eta gorago aipatu dugun taxonomia kontzeptuala eraikitzea lirateke helburuak. Horretaz gain, hitz-adieren desanbiguzioa, hiztegi aruntetarikoko informazio-eraztea, ezagutza lexikalaren errepresentazio-moduen azterketa, eta giza-erabilerarako hiztegi adimendunak eta itzulpen lexikaleraiko laguntza-sistemak.

9 Egungo Euskararen Bilketa Sistematika.