

Zer entzun, hura idatzi

Kortabitarte Egiguren, Irati
Elhuyar Zientziaren Komunikazioa

Idatziz jasotakoa bilatzea erraza da sarean. Horretarako, kontsultatu nahi dugun hitza bilatzailean idaztea besterik ez dugu. Bilaketa horietan, ordea, audio-fitxategietan esandakoak galtzen ditugu besteak beste, betiere, audio-fitxategi horietan esandakoaren azalpenak testu idatzian jasotzen ez badira.

AHOZKO HIZKETA EZAGUTZEA ETA HURA TESTU BIHURTzea EZ DA LAN ERRAZA. Hitzak ez dira ongi bereizten bata bestetik, intonazioa kontuan izan behar da, eta, gainera, seinale fisikoen zarata ere oztupo da. Horren harira, merkatu handia zabaldu da ahozko hizketa prozesatzen eta ulertzen duten sistementzat. Alegia, ahozkoa testu idatzi bihurtuko diguten tresnentzat.

Sistema horiek batez ere telefono bidezko zerbitzuetan integratzen dira oraingoz: aurretiko hitzordua, produktu-eskaerak, ikuskizunetarako erreserba-eskea eta abar. Baina badaude beste-lakoak ere: diktaketa automatikoa, adibidez. Azken horretan dihardute lanean, hain zuzen ere, EHUko Sistemen Ingeniaritza eta Automatika sailan, besteak beste.



ETBko *Gaur Egun* programak erabiltzen dituzte, besteak beste, hizketaren tratamendua egiteko sistemak trebatzeko.

Hizketaren tratamendua egiteko sistema asko eta ongi trebatu behar da. Alegia, sistemak nolabaiteko entrenamendua jaso behar du, makina-ikas-keta deritzona. Horretarako, batetik, telebista nahiz irratietako fitxategiak, audioak nahiz soinuak behar dira; eta bestetik, komunikabide horietan esan denaren erreferentziako testuak. EHUko ikertzaileek, adibidez, ETBko *Gaur Egun* eta *Teleberri* programak erabiltzen dituzte maiz, sistema trebatzeko. Ez da beharrezkoa hitzez hitz zer esan den jakitea; bai, ordea, esandakoaren laburpen bat jasotzeko gai izatea sistema. Azken finean, soinu eta hitzen arteko erlazioa ulertzen saiatzen da.

Ikasketa-prozesua amaitu ostean, edozein *Gaur Egun*-etan edo *Teleberri*-tan esandakoa ulertzeko gai behar luke izan sistemak. Ikastea prozesu motela izan arren, sistemak behin arauak edo informazioa barneratuta duenean, hau da, erreferentziako material egokia duenean, nahiko azkar erakusten du emaitza. Kasu honetan, ahoz esandakoaren testu idatzia. Azken finean, helburua da audio edo soinu batetik testua lortzea.

Txikia handi

Egia da merkatuan aurki daitezkeen horrelako aplikazio gehienek hizkuntza



Proiektua

Proiektuaren laburpena

Ikerkuntza-talde hau hizketaren ezagutza eleaniztunen alorrean aritzen da, euskararako eta haren inguruko hizkuntzetarako. Berezi, euskal komunikabideetako albistegien hizketatik informazioa automatikoki eskuratzeko hainbat tresna eta baliabide garatzen dituzte. Horretarako, informazio hori ahalik eta modu eraginkorrenean eskuratzeko teknikak ikertzen dituzte, eta, batez ere, baliabide urriko hizkuntzetarako metodoak garatzen dituzte —hala nola euskararako—.

Zuzendaria

Miren Karmele López de Ipiña doktorea.

Lantaldea

M.K. López de Ipiña¹, N. Barroso¹, N. Gilisagasti¹, I. Ariztimuño¹, A. Ezeiza¹, N. Ezeiza² eta M. Hernández².

Saila

Sistemen Ingeniaritza eta Automatika.

Fakultatea

¹Donostiako Unibertsitate Eskola Politeknikoa eta
²Informatika Fakultatea.



Taldea



Ezkerretik hasita, Ixabel Ariztimuño, Nora Barroso, Aitzol Ezeiza, Karmele Lopez de Ipiña eta Nerea Ezeiza.

'handiak' dituztela helburu; ingelesa, batik bat. Dena den, Donostiako Unibertsitate Eskola Politeknikoko iker-tzaileek, EHUko IXA, GTTS eta Adimen Konputazionala taldeekin elkarlanean, euskararekin dihardute lanean. Hizkuntza 'handi' eta 'txiki' horien arteko ageriko ezberdintasun nagusia erreferentziako datu-kopuruan datza. Mota horretako ingelesezko tresnek ikaragarriko datu piloa izaten dute; euskarazkoen erreferentziako materiala, berriz, dezente txikiagoa da. Horregatik, datu gutxi horiek hobeto eta zehaztasun handiagoz aprobetxatzeko teknika berriak bilatzen ari dira iker-tzaileak.

Zehaztasun-maila hori lortzeko, zenbait ekuazio matematiko erabiltzen dituzte. Datu-multzo eta audio-fitxa-

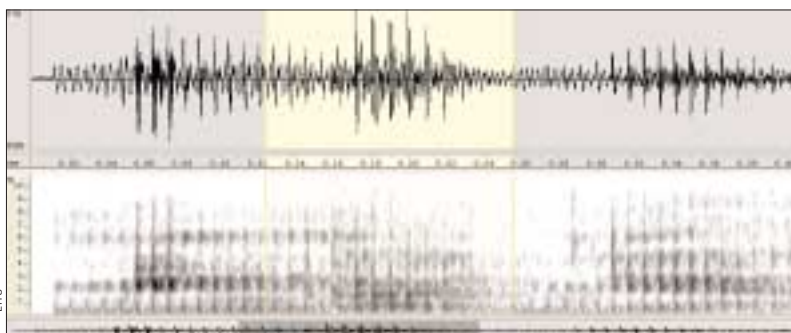
“helburua da audio edo soinu batetik testua lortzea; alegia, soinuen eta hitzen arteko erlazioa lortzea”

tegietatik informazio aproposa eman-go duten ezaugarri garrantzitsuenak aurkitzen saiatzen dira. Dena den, nahiko zaila da hautaketa hori egitea; alegia, jasoko den eta baztertuko den informazioa aukeratzea. Normalean, maiztasunarekin eta intonazioarekin lan egiten dute, une bakoitzean sis-

tema jasotzen ari den informazio-mota bereizteko (galdera bat edo adierazpen-perpau bat den bereizteko, adibidez).

Sistema horiek hizkuntzaren mende daude erabat, eta hizkuntza bakoitzak bere tresna du. Baina, EHUko ikertzaileek, adibidez, euskararekin ez ezik, gaztelaniarekin eta frantsesarekin ere egiten dute lan. *Teleberri* programak edo *Infozapir*-ko saiok aztertzen dituztenean, esaterako, bi helburu nagusi dituzte: batetik, gaztelania eta frantsesa ulertu nahi dituzte —euskararekin batera—, eta, bestetik, mota horietako sistemetan euskararen eta beste bi hizkuntza horien artean dauden antzekotasunak bilatu nahi dituzte, euskarazko tresnak hobeto trebatu ahal izateko.

Bide horretan, gaur egun, tresna beren hizkuntza bat baino gehiago erabiltzeko aukera aztertzen duten hainbat saiakuntza egiten ari dira. Horixe da, hain zuzen ere, EHUko ikertzaileen etorkizuneko erronka: euskara, gaztelania eta frantsesa ulertzeko gai izango den sistema bat garatzea. □



Ahoz esandakoaren maiztasuna eta intonazioa lagungarri dira sistema jasotzen ari den informazio-mota bereizteko.