

# Un corpus pal asturianu

Les tecnoloxíes llingüísticas na consolidación  
de les llingües minorizadas

*por* ROSER SAURÍ COLOMER

«The old quip attributed to Uriel Weinreich, that a language is a dialect with an army and a navy, is being replaced in these progressive days: a language is a dialect with a dictionary, grammar, parser and a multi-million-word corpus of texts — and they'd better all be computer tractable». (NICHOLAS OSTLER, *Foundation for Endangered Languages*)<sup>1</sup>.

## I. INTRODUCCIÓN

EN CUANTES que tresmisores d'una visión del mundu y de la historia, el poder de les llingües como aglutinantes sociales y expresión d'una identidá particular ye tan fuerte que les vuelve bien de veces nun problema pa los gobiernos d'imperios o entidaes polítiques multilllingües, que de cutio tomen midies p'amenorgar la diversidá: por exemplo, arrequexando la comunidá llingüística y culturalmente dixebrada nun apartaz de la realidá que se construí como país; cuándo creando

---

<sup>1</sup> Por cuenta de la reseña del *Workshop on Language Resources for European Minority Languages*, que tuvo llugar na *Internacional Conference on Language Resources and Evolution*, LREC 1998.

una tendencia a la marxinación social de los individuos y los pequeños grupos más marcadamente disidentes.

En primer casu, la llingua minorizada tuvo la posibilidá de sobrevivir hasta anguaño, apartada a un mundu fechu explícitamente de menos pola parte dominante y precisamente por esto pudo caltener la so estructura social. Fai'l casu de munches de les llingües indíxenes del continente americanu. En segundu casu, irónicamente la llingua pasa de ser un problema pal poder dominante a ser un problema pa los sos falantes: la continua penalización que reciben por usar la so llingua, dende los primeros años d'escolarización hasta los ámbitos profesionales y sociales de la vida adulta, fai qu'ésta dexa de ser el so mediu de comunicación natural y se vuelva una mena de desgracia de la qu'avergoñase, como ún d'esos alifaces familiares que-y toquen a ún ensin lo pidir. A diferencia del casu anterior, esta situación ye da-vezu propiciada cuando hai un nivel de contactu mayor ente les dos cultures.

La vitalidá d'una llingua ponse entós en peligru cuando se fai una torga pa los sos falantes. Esta ye, llamentablemente, la situación na que s'afayen anguaño munches de les llingües que se falen en mundu, incluyíes les qu'hasta agora sobrevivieren escaecíes en sociedaes consideraes arcaiques y ensin capacidá de futuru. Ello débese al contestu globalizador nel que nos alcontramos, qu'amista un nivel adicional de diglosia onde yá esistía ún previu porque, ente munches otres coses, impón un procesu d'estandardización cultural que pasa tamién pela creación d'una norma d'usu llingüísticu a favor de les llingües consideraes de más prestixu; esencialmente, les poques con capacidá de trescendencia a los sistemes d'información y con presencia nel desenvolvimientu tecnolóxicu occidental.

Sicasí, la mesma presea cola que s'establez esti procesu unificador ye la d'abrir la posibilidá de garantizar la supervivencia de les llingües minorizaes. La popularización de les tecnoloxíes de la informa-

ción y comunicación telemática abriéron-y camín al marcu actual de la llamada Sociedá de la Información. De la mesma, l' accesu a los medios de comunicación de mases, per un sitiu, y a les tecnoloxíes pal procesamientu automáticu del llinguaxe humanu, per otru, permitirán en bona medida la normalización social de les llingües minorizaes, en cuantes que les autorizarán como llingües afayadices pa situaciones comunicatives prestixaes. El desenrueldu de ferramientos y recursos pal procesamientu del llinguaxe ye, poro, un arma pa la supervivencia y la relevancia de les comunidaes llingüístiques minorizaes.

El presente artículu céntrase na rellación potencial de la llingua asturiana col segundu d'estos preseos: les tecnoloxíes llingüístiques. El puntu de partida que m'afala a escribilu ye l'enfotu en que l'asturianu, a diferencia d'abondes llingües del mundu con un númeru reducíu de falantes y en situación de subordinación respecto a otra, entá ta a tiempu d'entrar tamién nel dominiu de les tecnoloxíes llingüísticas.

Otres llingües minorizaes, xeográfica y políticamente averaes al asturianu, yá entraren nesti terrén: el gallegu, l'euskera y el catalán. Les condiciones de partida d'estes foron relativamente más favorables que les de la llingua asturiana, por mor del so estatus de co-oficialidá dientro del marcu políticu nel que s'alcuentra la mayor parte del so territoriu y fala. Sicasí, la esperiencia en toes elles apurre un bon puntu de referencia pa la entrada del asturianu a les tecnoloxíes del llinguaxe. Siguiendo estos trabayos previos a lo mesmo que la esperiencia algamada n'otres comunidaes llingüísticas, vemos que l'inquiz pa la entrada d'una llingua a estes tecnoloxíes ye la construcción d'un corpus llingüístico. Nesti artículu plantego, entós, la posibilidá d'un proyeetu talu pal asturianu.

## 2. CORPUS LLINGÜÍSTICOS

### *¿De qué falamos cuando falamos de corpus?*

Un corpus llingüísticu ye una colección de testos orales o escritos d'una llingua, esbillaos en cuenta d'unos criterios llingüísticos explícitos, pasaos a soporte electrónico y mínimamente procesaos, y que se manden como amuesa representativa d'esta pal atropamiento de datos empuestos al so estudiu sistemáticu<sup>2</sup>. Dientro de delles aries dedicaes al estudiu del llinguaxe, l'usu de corpus como recopilaciones de datos pa la xera científica nun ye una práctica nueva. Por exemplu, al sen más axustáu al so étimu llatín, na estaya de los estudios lliterarios dízse-y corpus a la totalidá de la obra d'un determináu autor. De la mesma, coñozse como *corpus* el conxuntu de contestos atropaos por lexicógrafos (bien de veces provenientes de la obra de los clásicos), col envís d'iguar y exemplificar les definiciones de los vocablos de la llingua en cuestión.

Estos usos de *corpus* estrémense, sicasí, de lo qu'anguaño se reconoz como l'oxetu de la Llingüística de Corpus. Un *corpus* llingüísticu estrémase d'un *corpus* lexicográficu tradicional pol fechu de qu'esti últimu suel recoyer namá amueses fragmentaes de la llingua, anquequier mui infrecuentes y anecdótiques, mientras que l'otru respe por una mayor representatividá. Coles mesmes, los corpora llingüísticos estrémense d'otros *corpus* de datos testuales (como los constituyíos pola obra d'un autor) pola so finalidá: los primeros fáense col envís específico de sirvir d'encontu p'análisis llingüísticos, y poro diséñense por aciu de criterios llingüísticos explícitos; mentes que los segundos entámense a partir de

---

<sup>2</sup> Esta definicion mira d'apautar les distintas caracterizaciones de corpus llingüísticos na bibliografía. Vease, por exemplu, SINCLAIR (1996), TOGNINI-BONELI (1996), TORRUELLA y LLISTERRI (1999) o McENERY y WILSON (2001).

criterios esternos a la llingua de los testos. Otramiente, los corpus llingüísticos caracteríicense pol tratamientu col que ye penerada la so información al envís de facilitar la xera d'ánalisis. Con éses, independientemente del soporte orixinal de los testos que los componen –oralidá o escritura– los corpus llingüísticos son informatizaos y, como se verá, na mayoría de los casos tamién procesaos mínimamente por facilitar l'accesu y la recuperación de la información.

El papel que xuega la dixitalización de la información nos corpus testuales nun ye menor. Cómo ye, que dellos autores refiérense a los corpus llingüísticos actuales como corpus informatizaos. Esto tien la so xustificación en contestu histórico de la disciplina. La utilización de colecciones de datos llingüísticos pal analís sistemáticu d'una llingua algamó reconocencia dientro de la tradición de la llingüística de campu a principios del sieglu, mesmo que na escuela estructuralista posterior. Sicasí, esta práctica féxose de menos col surdimientu de los planteamientos de xeitu mentalista defendíos por Chomsky dende finales de los 50, y que marquen un xiru nos estudios llingüísticos contemporanios. L'oxetivu d'esti nuevu paradigma ye l'estudiu de la gramática universal subxacente a la competencia de los falantes, y coles mesmes, l'actuación contémplase namá como'l reflexu imperfectu del conocimientu llingüístico. En cuantes que por fuerza reproducen esti nivel d'actuación, los conxuntos de datos llingüísticos atropaos na tradición estructuralista y de la llingüística de campu dexaron d'interesar.

Pero, magar la so impopularidá, los estudios de calter más empíricu nun se desdexen dafechu. Entemás que mientres les décadas de los 60 y 70 dieron n'iniciar, como elementu indisociable a esti tipu de planteamientu, la utilización de recursos tecnolóxicos. Con éstes, ente que nos trabayos previos de los estructuralistes y los llingüistes de campu'l soporte material de la información yera'l papel, darréu'l mediu d'almacenaxe vien ser electrónicu y la manipulación de los datos puede beneficiase de procesos automatizaos. Paralelamente, los

averamientos de base empirista vuelven pasu ente pasu a ganar adeptos, a cuenta de razones tanto de mena metodolóxica (por casu, aries como l'alquisición del llinguaxe nun son de sofitase n'intuiciones de llingüista), como de formulación teórica de base sobre la naturaleza del llinguaxe natural (una bona introducción a tou esti procesu ye McEnery y Wilson 2001). D'esta miente, nos años 80 el campu conoció anguaño como la Llingüística de Corpus espoxiga gracias a la converxencia d'esa perspectiva adoptada en dellos ámbitos de la llingüística coles posibilidaes de tratamiento de datos que permite'l desenrueldu tecnolóxicu del momentu.

### *Esbilla la información: clases de corpus*

En cuantes que los corpus llingüísticos tienen como finalidá apurrir información pal estudiu d'una o más llingües, mírase a que constituyan fragmentos representativos d'estes. La representatividá d'un corpus ye ún de los aspectos sobre los que más se tien aldericau na bibliografía del aria, al rodriu de consideraciones como: qué tamañu tien de tener un corpus que sía daveres representativu; qué clas de testos ha contener; en qué cantidá, etc. Véase d'exemplu Quirk (1992), Biber (1993) o Sinclair (1986).

Mientras l'espoxigue de la Llingüística de Corpus, nes décades de los 80 y 90, la representatividá plantegóse esencialmente en términos de tamañu (según más información, mayor representación de la llingua), y d'equilibriu ente los distintos usos de la llingua (variantes dialectales, oralidá versus escritura, diversidá de rexistros y xéneros, etc.). Anguaño, el tamañu sigue valorándose como un elementu importante a la de capturar el mayor número de fenómenos posibles d'una llingua. Sicasí, la cata d'una representación permediada de toles variedaes d'usu foi acutándose únicamente al criteriu pa la construcción de los llamaos corpus xenerales (los de reflexar la llingua común na totalidá de les variedaes y ámbitos). Con éstes, la idea de que les aplicaciones

finales determinen diseños diferentes de corpus foi cuayando de mou natural, invalidando los plantegamientos qu' imponén una serie de criterios estándares pal algame de la representatividá.

Veremos les aplicaciones potenciales d'un corpus llingüísticu na sección viniente. Pel momentu, ye la d'introducir les posibles clases de corpus, afitaes a cuenta de distintos parámetros complementarios ente ellos<sup>3</sup>:

- *Soporte orixinal de los datos:* tenemos corpus escritos o corpus orales, segúن tean constituyíos por testos escritos o orales.

- *Llingües que lu constituin:* pueden estremase corpus monollingües y multillingües. Estos últimos denómense corpus paralelos a tar constituyíos polos mesmos testos en caúna de les llingües representaes.

- *Variación llingüística que se mira a representar:* diatópica, diastrática y/o diacrónica. Dellos corpus representen más una de les variedaes de la llingua. A otra mano tán los corpora xenerales, que miren de representar varies d'elles, xeneralmente acutaes cronolóxicamente. Bona cuenta, nesti tipu de corpus ye mester asegurar l'equilibriu de representación ente les distintes variedaes.

- *Nivel d'especialización de los testos:* los corpus especializaos concéntrense en testos de determinaes aries d'especialidá. La finalidá suel ser mui específica, xeneralmente dientro del ámbitu de la terminoloxía y la terminografía.

- *Criterios d'esbilla de los datos:* estrémense corpus constituyíos a partir de testos enteros y los que contienen namá fragmentos (xeneralmente de llargor constante), que se conocen como corpus d'amueses. Anguaño, esta dixeبرا tira a escaecese por mor de los los problemas

<sup>3</sup> La mio caracterizacion ye mínima y ta malapenes destinada a cubrir de mou básicu la ralura de que pueda cadecer un lector non familiarizáu col tema. Pa una clasificación más refecha ye de vese TORRUELLA y LLISTERRI (1999).

que planteguen los corpus d'ampus pa la recuperación del contestu global d'usu de les espresiones llingüístiques.

- *Criterios d'actualización de los datos:* dacuando interesa la renovación periódica de los materiales d'un corpus. Esto ye sobremanera en corpus que respon por representar la llingua actual, nos que'l fragmenetu de testos más antiguos ye regularmente sustituyíu por otru de materiales recientes. Esti tipu de corpus conozse como corpus monitor.

### *Utilización de la información: aplicaciones posibles d'un corpus*

Les utilidaes qu'ufierten los corpus llingüísticos benefician a distintos campos y disciplines rellacionaes col estudiu del llinguaxe. Pueden xuntase en tres grandes niveles: dende les ciencies del llinguaxe, un corpus val de base pa estudios de plantegamientu teóricu. Por casu, d'abastecedor de contestos pa la validación de modelos de descripción de la llingua, pal análisis de la rellación ente dos llingües (a partir de corpus paralelos), o bien como indicador de la variación y les tendencies d'usu, d'interés n'estayes como la dialectoloxía y la sociollingüística, y qu'en casu de les llingües minorizaes puede emponese a determinaes decisiones de planificación llingüística. De la mesma, un corpus puede utilizase pal desenrueldu de productos de base llingüística. Fai'l casu de la ellaboración o meyora de gramátiques y diccionarios (esbilla del léxicu más frecuente, usu d'exemplos reales, distinción de sentíos a cuenta de los datos y non de la intuición del lexicógrafu, etc.); na creación de vocabularios especializaos; o na construcción de ferramentes pal aprendizaxe d'una llingua estranxera.

Nun segundu nivel, los corpus mándense tamién nel desenrueldu de tecnoloxía llingüística. Por exemplu, los corpus paralelos úsense pa entrenar sistemas de traducción automática, los corpus orales nel ámbitu de reconocencia y síntesis de fala, y los testuales pal entrenamiento de correctores ortográficos y gramaticales.

Finalmente, na estaya del procesamientu del llinguaxe natural los corpus utilíicense tamién de manera retro-alimentativa, na meyora de los preseos de procesamientu de corpus (sobremanera, pal entrenamiento de los de base estadística), talos como analizadores morfolóxicos y sintácticos, o etiquetadores semánticos. Vemos darréu cuáles son estos ferramientos y el so valir.

### *Tratamiento la información: corpus anotaos*

Otra manera, los corpus llingüísticos nun son solo esto: un conxuntu de testos caracterizaos a cuenta de determinaos criterios, que s'embalaguen en soporte electrónico, y que son utilizaos con una finalidá específica. Precisamente a cuenta de la so finalidá como fonte d'abastecimiento de datos, ye mester que la información que contién sia bona d'algamar, y néstes la dixitalización de testos nun ye abondo. Porro, na mayoría de casos los corpus llingüísticos presenten, amás del so conteníu testual orixinal, un nivel metalingüístico d'información qu'esplicita les características gramaticales de los distintos niveles de constituyentes: dende trazos morfolóxicos de los elementos léxicos hasta, potencialmente, características de los párrafos o unidaes testuales mayores, en casu de corpus escritos; nos corpus orales, dende información fonética hasta estructura prosódica. A cuenta de la mio experiencia personal dientro de la Llingüística de Corpus, nesta estaya y lo que sigue del artículu voi concentrarme casi dafechu na parte de corpus escritos. El llector, sicasí, afayará abondes referencies bibliográfiques a la d'enanchar la so conocencia tocantes a corpus orales. Véase por casu Llisterri (1997, 1999).

### *Niveles d'anotación*

Los corpus llingüísticos qu'apurren esti nivel de metainformación suelen referise como corpus etiquetaos, cuidao que la información co-

difícilase en testu pentemedies de códigos específicos o conxuntos d'etiquetas (o etiquetarios). De la mesma, tamién puede dicise corpus anotaos<sup>4</sup>. Per un sitiu, l'etiquetaxe d'un corpus facilita la xera de remanar y recuperar información. Por casu, pa un lexicógrafu o un gramáticu que quiera analizar los contestos d'usu d'un verbu determináu, disponer de toles formes d'esti verbu identificaes pola so rellación con un único lema quita de tener que buscar los contestos d'usu a partir de la enumeración refecha de toles formes posibles. Otramiente, l'etiquetaxe permite una cuantificación de datos más sofisticada que la pura frecuencia d'usu de cada forma léxica (por exemplu, distribución de tiempos verbales, clases de complementos pa determinaos predicaos, tipos de modificación aplicada a una clas de nomes particular, coocurrencias frecuentes d'elementos léxicos, variantes dialectales, etc.). Con éstes, pueden albidrarse tendencias de comportamientu llingüísticu dentro del fragmentu de llingua analizáu.

L'etiquetaxe d'un corpus suel informar sobre distintos niveles: l'etiquetaxe estructural, por casu, marca la organización estructural del testu: título, subtítulo, capítulo, párrafu, sección, subsección, etc..., hasta frase, xeneralmente. L'etiquetaxe morfolóxicu aplícase al nivel de los elementos léxicos, indicando la so categoría grammatical y, en llin gües como les romances, tamién la información de flexón verbal y nominal. Xeneralmente, sobre esti nivel aplícase l'etiquetaxe sintácticu, que puede realizase con más o menos fondura, d'acordies cola aplicación pa la que se plantega'l corpus. Con éstes, l'anotación puede ser superficial, de simples grupos nominales y/o verbales, o más comple-

---

<sup>4</sup>Nos últimos años, l'usu de corpus non etiquetaos fíxose tamién una práctica relativamente avezada, pola economía de recursos y tiempu que conlleva. Sicasí, na mayoría de casos trátase de corpus utilizaos pal entrenamientu de ferramentas estadístiques de procesamientu del llinguaxe (veremos más alantre a qué me refiero), en vez de corpus pensaos pal desenrueldu de preseos de base.

xa, indicando les funciones sintáctiques de los distintos componentes d'una fras (suxetu, oxetu, modificadores, etc.)<sup>5</sup>.

Éstos son en grandes trases los niveles d'anotación más comunes. Por embargu, otros niveles d'información pueden faese igualmente prestosos d'acordies cola aplicación específica a la que s'empobine'l corpus. Por exemplu, puede introducise tamién anotación de calter semánticu o pragmáticu. Pero, según nos niveles anteriores podía trabayase dende un plantegamientu en forma neutru tocantes a cualquier teoría llingüística, dende'l nivel semánticu l'etiquetáu mira o a basase nuna visión particular de la semántica o la pragmática, o a especializase nun aspectu concretu, como les rellaciones anafóriques, les unidaes d'información temporal, los marcadores discursivos, los papeles temáticos de los predicaos verbales, etc.<sup>6</sup>.

### *Ferramientes de procesamientu de corpus*

Toos estos niveles d'etiquetaxe fáense al traviés de varies capes de procesamientu automatizáu. Por exemplu, pa l'anotación morfolóxica empléguese los programes conocíos como analizadores

<sup>5</sup> Una y bones en forma bibliografía d'esti campu vien n'inglés, cuento importante introducir equí la terminoloxía de los distintos niveles d'etiquetaxe tamién nesta llingua. A la fase d'anotación morfolóxica conózse-y como *part-of-speech* (POS) *tagging*. Al nivel d'etiquetaxe sintácticu, como *parsing*, identificándose tamién el marcase sintácticu superficial como *shallow parsing* o bien *chinking*.

<sup>6</sup> Exemplos particulares de corpus nesti grupu son: el trabayu presentáu en FLIGELSTONE (1992), que se caracteriza pola anotación de rellaciones anafóriques; el *RST Corpus* (CARLSON *et al.* 2003), nel que los testos anótense a partir de la teoría de discursu *Rethorical Structure Theory* (MANN y THOMPSON, 1988); el *TimeBank* (PUSTEJOVSKY *et al.* 2003), nel que namá s'anoten les eventividáes y les expresiones temporales col envís de mandase d'elles n'estudios sobre razonamiento temporal; o tamién el *Corpus de Operaciones Metalingüísticas Explícitas* (RODRÍGUEZ 2003), nel que s'anotaron les operaciones de creación de conocimientu nel marcu del discursu científicu.

morfolóxicos, que trabayen cola collaboración d'un diccionariu en formatu electrónico nel que cada elementu léxicu xúncese a una o más categoríes gramaticales y, en llingües como les romances, d'un lematizador, que ye'l d'encargase de dixebrar nes palabras la raíz y les terminaciones. Na xera d'etiquetaxe sintácticu úsense los llamasos analizadores sintácticos, que como yá se comentó, pueden trabayar a un nivel más o menos fondu d'análisis: dende la simple detección de grupos nominales y verbales, hasta'l marcase de funciones y rellaciones de dependencia. Pa l'anotación del nivel semánticu (y mesmo'l pragmático), nun hai ferramientes específiques que permitan l'agrupación bajo un único denominador, según asocedía co los niveles previos. Sicasí, vaga la pena comentar la esistencia de los etiquetadores semánticos, sobre los que se tien trabaya en forma apocayá. Trátase de ferramientes qu'asignen etiquetas semántiques a les unidaes léxiques a partir de la so clasificación nuna ontoloxía determinada, lo que, naturalmente, quier un llabor previu de desambiguación léxica.

D'últimes, ye importante mentar tamién les ferramientes pa la esplotación de corpus, que son básiques na investigación lexicográfica y llingüística. Les más destacaes equí son les que realicen busques de (secuencias de) palabras, presentando los resultaos ordenaos alfabéticamente y col contextu oracional, y ufiertando arriendes la posibilidá de cómputu de frecuencies.

Les ferramientes de procesamientu de corpus estrémense en xeneral según el plantegamientu de base de que partan: d'uno, tán les ferramientes de base simbólica; d'otru, les de base estadística. Si ye les primeres, desenvuélvense a partir de conocimientu llingüístico, y poro, tan en contactu con delles de les vertientes de la llingüística teórica. Les segundes básense en modelizaciones estocástiques del fragmentu de llingua col que se trabaya. El conocimientu llingüístico utilízau nesti segundu casu ye mínimu, si non inexistente, y poro nun

dexa, desque desenroldada l'aplicación en cuestión, faer una abstracción en términos llingüísticos de los fenómenos trataos. Sicasí, suel dar resultaos meyores que'l d'una ferramienta de base simbólica desenroldada per un períodu de tiempu equivalente.

Arriendes del averamiento simbólicu o estadísticu al problema, el campu dedicáu al desenvolvimientu de toes estes ferramientos pal tratamientu de testos escritos conozse como Procesamientu del Llinguaže Natural (PLN) o, en dellos contestos tamién, Inxeniería Llingüística; mentanto que no que fai al procesamientu de testos orales fálase de Tratamientu de la Fala. Equí voi emplegar el térmigu xenéricu de tecnoloxíes del llinguaxe (o'l ximielgu de tecnoloxíes llingüístiques), a la de referime globalmente a estos dos aries en xunto a la Llingüística de Corpus. Al dicir del so nome, el so denominador común ye l'usu de tecnoloxía pal tratamientu del llinguaxe.

### 3. ESPERIENCIAS DE REFERENCIA

Desque afitaos los preliminares, esta estaya presenta en grandes tráces les esperiencias previes de proyectos de corpus que cuido de más interés pa lo qu'equí nos lleva'l tiempu: la construcción d'un corpus llingüístico pa la llingua asturiana.

Les primeres esperiencias que derrompen la nueva dómina na Llingüística de Corpus danse pa la llingua inglesa. Anguaño, constitúin obres de referencia y, magar que tean bonamente documentaes na bibliografía de la disciplina, vaga bien presentales mínimamente. La so naturaleza de finxu ye, claramente, resultanza del calter pioneru de la investigación en delles comunidaes onde se fala esta llingua. Esisten, bona cuenta, corpus n'otres llingües, cuándo cercanes xeográficamente (como'l francés o l'alemán), cuándo más alloñaes, dende'l xaponés, el coreanu y el mandarín, hasta les recientemente tan preciaes llingües falaes en mundu musulmán, como l'árabe y el

farsi<sup>7</sup>. Cuento, sicasí, que queda fora del algame d'esti trabayu faer un repás, anquequier mínimu, de los casos más relevantes pa cada llingua.

Otra manera, voi revisar tamién los proyectos esistentes dentro del marcu xeográficu de la Península Ibérica, centrándome dafechu naquellos que se desenvuelven al abellu del mesmu contestu políticu nel que s'inclúi l'asturianu; esto ye, l'Estáu español.

### *Los pioneros*

Arriendes de la mesma evolución del campu, nos corpus del inglés son d'estremase dos xeneraciones. La primera inclúi los corpus entamaos nes décades de los 50 y 60, de la que los averamientos empíricos al llinguaxe dexaron d'interesar, particularmente dientro de la llin güística desenrolada nel contestu de fala inglesa. Los corpus d'esti períodu inicial caracterízense por ser de tamañu modestu en comparanza con dellos de los corpus construyíos de más recién. Los más destacables son:

- *Survey of English Usage*: Corpus d'inglés británico entamáu en 1955 y desenrolzáu a lo llargo de, al aldu, 30 años, nel University College London. Contién un millón de palabras, ente testos orales y escritos. Ye de remarcarse que los materiales d'esti corpus emplegáronse como base pa la constitución de la conocida gramática de referencia del inglés (Quirk *et al.* 1985).

---

<sup>7</sup>A esti sen, vemos que sigue siendo parte de la comunidá de fala inglesa la de decidir promocionar recursos nuna o otra llingua. Ye indicativo la creciente demanda de llingüistes, traductores ya intérpretes en delles d'estes últimes llingües xenerada de magar el 11 de setiembre del 2001. L'emburrión a la investigación en delles llingües vese afaláu por intereses de la comunidá de les axencies d'espionage ya intelixencia militar, y con oxetivos que poco o nada tienen que ver cola investigación puramente científica o l'apreciu de valores culturales.

• *Brown Corpus, Lancaster-Oslo/Bergen Corpus (LOB) y Kolhapur Corpus:* Corpus de testos escritos nes variantes d'inglés americano, británico ya inglés de la India, respectivamente. El primeru atropóse na década de los 60 y los otros dos n'etapes posteriores, como corpus equivalentes al anterior pa otres variantes del inglés. Contienen toos un millón de palabras y presenten los testos agrupaos en 15 categoríes diferentes, procurando d'esta miente una bona representatividá de los distintos usos de llinguaxe escrito. Son, poro, corpus con vocación de referencia (Francis y Kucera, 1964).

• *London-Lund Corpus:* Corpus oral d'inglés británico, de 500.000 palabras. Contién esclusivamente trescripciones de material falao, orixinario de dos proyectos diferentes: de la parte oral del Survey of English Usage presentáu enantes, y del Survey of Spoken English, entamáu na Lund University en 1975 como proyeetu hermanu del anterior (Svartvik 1990).

Anguaño, la dimensión de los corpus tien variao enforma por mor de los avances tecnolóxicos. La capacidá d'atropamientu ye bramente mayor de lo que yera de la qu'empezara la tecnoloxía dixital, y los procesos de tratamientu de datos son abondo más rápidos. A quasi 50 años de magar la creación del primer corpus modernu, los corpus qu'hai disponibles pal inglés son muchos y d'estremada mena. D'ente ellos interesa mentar los vinientes pol so calter de referencia dientro del aria y la so trescendencia na estaya de la llingüística aplicada:

• *Bank of English:* Entamáu en 1991 de parte de la Birmingham University y COBUILD, una división de la editorial Harper Collins. El so principal envís ye mandase como fonte de datos n'estudios llingüísticos y la fechura de diccionarios. Trátase d'un corpus monitor; esto ye, un corpus al que se-y amiesta dacuando material nuevo. La so última edición, del 2002, atropa 450 millones de palabras (Sinclair 1987)<sup>8</sup>.

---

<sup>8</sup>Véase tamién [http://titania.cobuild.collins.co.uk/boe\\_info.html](http://titania.cobuild.collins.co.uk/boe_info.html).

• *British National Corpus* (BNC): Contién un total de 100 millones de palabras del inglés moderno, tanto del rexistru escritu como del oral, qu'inclúi hasta parolaes coloquiales. Foi atropáu de parte d'una xermandía formada por editores británicos y instituciones académiques como la Oxford University y la Lancaster University. Tanto nesti casu como nel anterior, el so diseñu respe pola mayor representatividá posible (Aston y Burnard 1998)<sup>9</sup>.

### *La Llingüística de Corpus nel nuesu ámbitu xeográficu*

El desenrueldu de corpus llingüísticos, mesmo qu'otres ferramentes pal tratamientu del llinguaxe natural nel ámbitu del Estáu Español, reproduz la división ente llingües oficiales o co-oficiales, d'uno, y llingües ensin reconocencia oficial, d'otro. En mayor o menor grau, el primer grupu de llingües (euskeru, catalán, español y gallegu) dispón anguaño de presea y recursos de procesamientu del llinguaxe, grupos d'investigación más o menos afitaos aplicaos al so desenvolvimientu, y (non por último, menos importante) la voluntá institucional de sofitar esta llinia de trabayu. Sicasí, la Llingüística de Corpus y el Procesamientu del Llinguaxe Natural son tierres entá non derrotes en casu de llingües ensin oficialidá como l'asturianu.

El Procesamientu del Llinguaxe Natural nel nuesu ámbito surde metá los 80, al rodriu de los proyectos de traducción automática que s'entamen en Barcelona y Madrid, tanto dende empreses privaes como financiaos públicamente pola Comisión Europea. El campu pasa daquella per una dómina d'euforia paralela a la espectación qu'esiste nel marcu européu más xeneral, pero, desque nos 90, esos proyeutos dan n'abandonase porque la probeza de los resultaos obtenidos hasta

---

<sup>9</sup> Véase tamién: <http://www.natcorp.ox.ac.uk>.

entós, en comparanza colos recursos invertíos, apunta un fracasu nos plantegamientos de base<sup>10</sup>.

Sicasí, de la qu'estos proyectos se desmantelen, yá hai la infraestructura creada en términos d'instituciones y recursos humanos a la de seguir la investigación empobinada agora a la constitución de recursos llingüísticos esenciales (como corpus, diccionarios en formatu electrónicu y bases de datos léxicos), y a la igua de ferramientes básiques pal procesamientu del llinguaxe. A esto tien que se-y amestar el fechu de que, al par del desendolcu de proyectos de traducción automática na década de los 80, yá esiste dalguna institución trabayando na estaya de la Llingüística de Corpus, como fa'il casu del proyectu del *Diccionari del Català Contemporani* col so *Corpus Textual Informatitzat* (descritu darréu). Coles mesmes, ye metanos los 90 de la que la Llingüística de Corpus espoxiga nel nuesu ámbitu xeográficu, cuayando tanto en contestu español y en Cataluña (onde surdieren los primeros enteinos en traducción automática), como en País Vascu y Galicia, col desenrueldu consiguiente de recursos para caúna de les llingües falaes nestes zones.

Ufierto darréu una amuesa representativa, anque non refecha, de los corpus llingüísticos esistentes anguaño d'estes cuatro llingües: euskera, catalán, español y gallegu. Bien que nenguna d'elles presenta un panorama asemeyáu al del inglés, tampoco nun se caractericen pola situación llaceriosa na que s'afayen les llingües ensin reconocencia oficial. De la mesma, albídrase daqué disparidá nel nivel de desenrueldu del aria en caúna de les llingües. Ello correspuende en llinies xenerales colos datos que s'apurren na parte dedicada al fomentu y desenvolvimientu del aria dientro'l plan nacional d'investigación científica del gobiernu vascu<sup>11</sup>. Con

<sup>10</sup> Véase ABAITUA (1999) pa una refecha caracterización d'esti períodu.

<sup>11</sup> *Plan Nacional de Investigación Científica, Desarrollo e Innovación (2000-2003)*. «Propuesta de acción estratégica en el área sectorial “Sociedad de la Información”: Industria de

encontu nun estudiu recién ellaboráu de parte'l Consorciu Européu EURONMAP, ésti detalla qu'el repartu d'esfuerzu d'investigación y desenvolvimientu empobináu en caúna de les llingües en cuestión en marcu conxuntu del estáu ye ésti: 60 % pal español, 20 % pal catalán, 8,5 % pal gallegu y finalmente 7 % pal euskera. Sigo esti orde de llingua de más a menos recursos na presentación de los proyectos de corpus en caúna d'elles. Sicasí, voi centrame más nos proyectos desenrolaos pa les llingües co-oficiales que nos del español: anque en menor grau, parten col asturianu la so situación minoritaria y minorizada y, poro, la so esperiencia dentro de la Llingüística de Corpus resulta de gran utilidá a la d'apuntar un proyectu nesti terrén, con recursos y sofitu institucional acutáu.

### *Corpus d'español*

Esisten dellos corpus pa esta llingua. Dos de los más destacaos son los corpus de referencia atropaos pol Instituto de Lexicografía de la Real Academia Española, ún diacrónico y otru del español actual. Dos referencies que detallen el desenrueldu d'estos proyectos son Pino *et al.* (1999) y Sánchez-León *et al.* (1999). Dambos corpus inclúin documentos escritos en toles variantes, tanto peninsulares como estrape-nisulares, y tán eminentemente empobinaos a mandase como recurso básico nel trabayu lexicográficu d'esta institución:

- *Corpus Diacrónico del Español (CORDE)*: Corpus que mira de representar la variación del español a lo llargo de la so historia escrita. Contién documentos de tres dómines básiques: Edá Media, Sieglos d'Oro y Época Contemporánea (de mano, hasta 1975). Anguaño presenta un total de 145 millones de palabras y puede consultase per internet (<http://www.rae.es/cordenet.html>).

---

la Lengua e Ingeniería Lingüística». Eusko Jaurlaritza (Gobierno Vasco). Descargable en: <http://www.euskadi.net/euskara/datos/azkeninforme.pdf> (última actualización: 26/09/2003).

• *Corpus de Referencia del Español Actual (CREA)*: Corpus monitor, constituyíu de testos representativos de les distintes variantes del español actual. Comprende un trechu de 25 años, que s'actualiza de xemes con materiales más recientes. Esto significa que los documentos más vieyos treslládense al corpus CORDE. Anguaño contién 145 millones de palabras.

Esisten otros corpus del español entamaos por editoriales, como por exemplu:

• *Corpus CUMBRE*: Desenrolzáu pola Editorial SGEL con fines lexicográficos. Atropa 20 millones de palabras, vinientes de testos tanto escritos como orales del español peninsular ya hispanoamericano. Los testos escritos pertenecen a distintos xéneros (lliterariu, periodísticu, ensayu) y toquen distintes estayes d'especialidá (filosofía, historia, ciencia, derechu, economía, etc.). Los testos orales proceden de grabaciones de programes de radio y televisión. Según les amueses orales son de la década de los 90, les escrites algamen el períodu ente los 50 y los 90 (Sánchez *et al.* 1995).

De la mesma, Vox Bibliograf trabaya col so corpus d'al rodíu de 10 millones de palabras, y la editorial SM con otru de 60.000 palabras. Amás d'estos corpus, esisten otros de menores dimensiones, cuando de llingua xeneral pero con un envís específico (ye la de los corpus d'editoriales que vienen de mentase), cuando de llinguaxe acutáu, como'l *Corpus de vocabulario del niño de 6 a 14 años* de la Universidá de Granada, o bien el *Corpus 92, Lengua escrita por aspirantes a estudios universitarios* de la Universitat Pompeu Fabra. D'últimes, vaga tamién mencionar la esistencia de corpus puramente orales, creaos pal desendolcu d'aplicaciones na estaya de les Tecnoloxíes de la Fala. Como yá anuncié enantes, nun voi centrarme nellos.

*Corpus de catalán*

Los corpus de mayor envergadura detállense darréu. Información adicional sobre corpus testuales nesta llingua puede afayase en Soler i Bou (1998) y, a nivel más xeneral, Llisterri (2000) ufierta una introducción abondo completa sobre los recursos y ferramientos de procedimientu del llinguaxe esistentes anguaño.

• *Corpus Textual Informatitzat de la Llengua Catalana (CTILC)*: Corpus de testos escritos de rexistru variáu (lliterariu y non lliterariu) desendolcáu nel Institut d'Estudis Catalans como componente básicu na creación del *Diccionari del Català Contemporani (DCC)*. Trátase, poro, d'un corpus de llingua xeneral, que respe pola representatividá y que foi diseñáu col envís de mandase como base na investigación lexicográfica. Compónenlu al aldu de 4000 testos dataos ente 1832 (fecha simbólica del aniciu de la dómina conocida como la *Renai-xença* de la llingua catalana, cola publicación de la *Oda a la pàtria*, de Bonaventura Carles Aribau) y 1988. En total, recueye 52,3 millones de palabras. El corpus ta lematizado y etiquetáu morfosintácticamente (Rafel 1992-93, 1996). Amás, compleméntase con una base de datos lexicográfica constituyida a partir de 13 diccionarios espublizaos en mesmu períodu, y sirvió yá pa la creación del *Diccionari de Freqüències* (Rafel i Fontanals 1996-1998).

• *Corpus del Català Contemporani de la Universitat de Barcelona (CUB)*: Corpus de testos escritos y orales, estructuráu en siete subcorpus independientes que se rellacionen ente ellos pol envís de configurar conxuntamente una caracterización del catalán actual dende la so variación diatópica, diastrática y diafásica. Contién: un subcorpus de variedaes xeográfiques (*Corpus Oral Dialectal*); ún de variedaes sociales (*Corpus Oral Social*); y los cinco restantes de variedaes funcionales (*Corpus Oral de Conversa Col·loquial*, *Corpus Oral de Re-*

gistres, *Corpus Oral de Publicitat*, *Corpus d'Informatius Orals*, *Corpus Escrit de Català Actual*). Los subcorpus escritos tán lematizados y anotaos morfosintácticamente. Los subcorpus orales tán transcritos y, según la naturaleza de los datos, presenten anotación discursiva, de grupos tonales, o bien lematización y etiquetáu morfosintácticu. Si en casu del CTILC, la utilidá que se quier ye la investigación lexicográfica, nésti mírase a crear una base pa estudios específicos sobre la variación del catalán. Bones introducciones al proyectu son Boix (1996) y Viaplana (2000).

• *Corpus textual plurilingüe especialitzat*: Corpus ellaboráu nel Institut Universitari de Lingüística Aplicada de la Universitat Pompeu Fabra (Bach *et al.* 1997). Contién testos escritos en catalán, español, francés, inglés y alemán, pertenecientes a les siguientes aries d'especialidá: mediu ambiente, informática, medicina, derechu y economía. La parte del catalán cuenta con unos 5 millones de palabras, lematizaes y con etiquetáu morfosintácticu.

Otros proyectos de corpus del catalán son: el *Corpus Parole*, de 21 millones de palabras desendolcáu en marcu d'un proyectu européu, y el *Corpus de Diatopia Diacrònica de la Llengua Catalana*, afaláu dende la Universitat de Barcelona. De la mesma, esisten dellos corpus de fala atropaos específicamente pal desenvolvimientu d'applicaciones de tecnoloxíes de la fala.

### *Corpus de gallegu*

Una panorámica xeneral de la situación del gallego ufiértala García-Mateo *et al.* (1998). Equí interesa destacar los corpus siguientes, todos ellos en tenores:

• *Corpus de Referencia do Galego Actual* (CORGA): Esti corpus ye ún de los proyectos entamaos en Centro Ramón Piñeiro para a Investigación en Humanidades, institución creada de parte de la Xunta de

Galicia por afalar la enseñanza, la investigación y l'usu de la llingua gallega. El CORGA plantégase como un corpus de referencia de testos escritos y orales, producidos ente 1975 y 2004 (fecha prevista pa la so finalización), y mira d'abrir diferentes rexistros del gallego actual. Anguaño embalaga unos 17,5 millones de palabras anque aspira a un total de 25 millones, anotaes cola categoría morfosintáctica. Los detalles del proyecto, mesmo que la última versión de CORGA, emitida en xineru del 2003, pueden consultase per internet<sup>12</sup>.

• *Corpus do galego moderno* (CGM): Corpus atropáu nel Instituto da Lingua Galega. Contién al rodiu de 10 millones de palabras, viñientes de testos escritos dende'l sieglu XVII hasta güei. Mira d'incluirse ende testos de la lliteratura oral, histories, refranes, canciones populares, etc. L'envís principal del corpus ye lexicográfico y pal estudiu de la llingua gallega en xeneral.

• *Arquivo do Galego Oral*: Corpus oral, igualmente atropáu pol Instituto da Lingua Galega. Estrémase en dos subcorpus: l' *Arquivo do galego popular* y l'*Arquivo do galego culto*. El primeru contién más de 1000 hores de grabaciones orales efectuadas ente 1974 y 1998, ya inclúi falantes de distintes edaes y variedades dialectales. El segundu inclúi grabaciones de charles, conferencies, y meses redondes sobre temes del ámbitu políticu y social.

• *Corpus Linguístico da Universidade de Vigo* (CLUVI): El CLUVI desendólcase en Seminário de Linguística Informática de la Universidade de Vigo (Aguirre Moreno *et al.* 2001, 2002, 2003). Trátase d'un corpus escrito y oral que recueye testos contemporáneos de distintos rexistros especializaos: xurídico-alministrativu, periodísticu, informáticu y lliterariu. Otra manera, los testos pueden ser monollingües en gallegu, traducciones gallegu-español y traducciones inglés-galle-

---

<sup>12</sup> <http://corpus.cirp.es/corga>.

gu. La parte de testos d'aniciu escritu estrémase en cuatro subcorpus d'al rodiu d'un millón de palabras caún: el corpus paralelu TEC-TRA, constituyíu de testos de rexistru lliterariu inglés-gallegu; el corpus paralelu LEGA, de testos xurídico-administrativos gallegu-español; el corpus monollingüe XIGA, de testos del ámbito informático escritos en gallegu; y finalmente el corpus monollingüe MEGA, del ámbito de los medios de comunicación social. De la mesma, trábáyase na construcción d'otra parte adicional constituyida de testos paralelos portugués-gallegu. Quier llegase al nivel d'anotación morfosintáctica, mesmo qu'a la alliniación oracional de los subcorpus paralelos. Les aplicaciones que se quieren construir en base a esti corpus son varíes: dende la construcción de ferramentas básiques como extractores d'información léxica, terminoloxica y fraseoloxica, hasta aplicaciones más sofisticadas y de llargu plazu, como la extracción d'información, la sumarización o la recuperación d'información *on-line*.

Ye de remarcarse tamién la esistencia d'otru corpus oral, VOGA-TEL, desenrolzáu de parte de la empresa Telefónica I+D en collaboración cola Universidade de Vigo. El so envís ye puramente l'entrenamiento de ferramentas pal tratamientu de la fala.

### *Corpus d'euskera*

L'euskera tien dos corpus llingüísticos principales, dambos de testos escritos. Preséntense n'Agirre *et al.* (1998), dientro del marcu xeneral pal procesamiento del euskera que se desenvuelve na Euskal Herriko Unibertsitatea. Esisten amás otros corpus de calter menor, confeccionaos de parte de grupos d'investigación por entrenar les ferramientes de base estadística que constrúin. Equí voi centrame má nos corpus de referencia, que son:

- *Orotariko Euskal Hiztegia* (OEH, ‘Diccionariu Xeneral Vascu’):

Corpus testual diacrónico que se manda como base na creación del diccionariu descriptivu d'Euskaltzaindia, l'academia de la llingua vasca. Atropa un total de 5.800.000 palabres, vinientes de 310 obres escrites en distintes variedaes del euskera, dende'l sieglu XVI hasta la d'entamase la estandardización de la llingua, contra 1960.

• XX. *Mendeko Euskararen Corpus Estatistikoa* ('Corpus Estadístico del Euskera del sieglu xx'): Corpus del euskera actual, constituyíu por amueses de testos vinientes d'al rodíu de 6000 obres escrites mientres el sieglu XX. En ficies de recoyer la máxima variedá léxica y estructural posible, establecióse en base a una esbilla aleatoria del inventariu d'obres escrites n'euskera, atropáu por UZEI (*Terminología eta Lexikografiako Zentroa, Centru Vascu de Terminoloxía y Lexicografía*). Estes obres pertenecen a toa triba de xéneros, dende poesía, teatro y lliteratura infantil, hasta investigación o llibros de testu. Amás, miróse a representar caún de los cuatro períodos relevantes de la historia del euskera del sieglu XX: 1900-1939 (dende l'entamu de sieglu hasta les guerres), 1940-1968 (de magar la posguerra a la nacencia del euskera estándar), 1969-1990 (dende los cambios producidos pol estándar hasta les publicaciones de les normes de Euskaltzaindia), 1991-1999 (posterior a la normativa). El corpus atropa un total de 4,7 millones de palabras, lematizaes col lema estándar y el de la variante dialectal, y etiquetaes cola categoría gramatical y les distintes acepciones. Diseñóse col envís de representar les distintes variedaes diatópicas y diafásicas del euskera contemporáneo. Poro, ta pensáu como fonte de datos qu'espeye l'usu real de la llingua, tornándose d'una función puramente prescriptiva<sup>13</sup>.

---

<sup>13</sup> Na páxina d'UZEI ([http://www.uzei.com/default\\_cas.html](http://www.uzei.com/default_cas.html)) ufiértase más información. Otra manera, el corpus ye consultable per internet ([http://www.euskaracorpusa.net-IXXmenda/Konts\\_arrunta\\_fr.html](http://www.euskaracorpusa.net-IXXmenda/Konts_arrunta_fr.html)).

#### 4. UN PROYECTU DE CORPUS PAL ASTURIANU

##### *Tecnoloxía llingüística y llingües minoritaries*

Un proyector de corpus llingüísticu pal asturianu tien de char cuenta necesariamente del so estatus de llingua minoritaria. En cierto, les posibilidaes de desendolcu de tecnoloxía llingüística pa una llingua de complexón sana son bien superiores a les que tien una llingua como l'asturianu o l'euskera. El problema básicu de les llingües identificaes como minoritaries ye qu'el nivel de recursos de que disponen, materiales y humanos, ye enforma menor qu'el de les mayoritaries. A ello tien que se-y amestar el ruin, si non nulu, interés comercial que presenten, lo que supón un costu mui altu pa la investigación y el desenvolvimientu nel ámbitu.

El problema nun ye, sicasí, esclusivu de les llingües minoritaries. Hai llingües mayoritaries que s'afayen na mesma situación, como por casu l'hindi, magar los sos 180 millones de falantes (Somers 2001). Amás, la etiqueta de llingua minoritaria presenta della ambigüedá. Si por minoritaria entendemos llingua de pocos falantes, debemos incluir llingües como'l finladés (4,7 millones en Finlandia) o'l suecu (7,8 millones en Suecia)<sup>14</sup>. Nestos casos, por embargu, la condición minoritaria nun quita un bon allugamientu nel ranking de llingües pa les qu'esisten aplicaciones pal procesamiento del llinguaxe: el finlandés alcuéntrase na sesta posición y el suecu na séptima, darréu del inglés, l'alemán, el francés, l'español y l'italiano, d'acordies el Natural Language Resource Registry (NLRS)<sup>15</sup>. Per otru sitiu, si de llingua minori-

<sup>14</sup> Datos vinientes de: <http://www.ethnologue.com>.

<sup>15</sup> <http://registry.dfkki.de>. AGIRRE *et al.* (2002) ufierta un análisis d'estos datos, empobináu a les aplicaciones de PLN independientes de llinguaxe.

taria entendemos llingua en situación minoritaria dentro del marcu políticu nel que s'afaya, el términu yera d'aplicase tamién a les llingües d'inmigrantes con rellación al so país d'adopción, anque sían de llingües mayoritaries nel so país d'aniciu.

Ye mester entós considerar cuálos son los elementos que garanticen o torguen l'accusu d'una llingua a les tecnoloxíes del llinguaxe. Veremos darréu cómo l'estatus de cada llingua determina de mou xeneral la rellación qu'ésta caltién coles tecnoloxíes del llinguaxe. Finalmente, identificaremos na sección siguiente les traces que caractericen la llingua asturiana y aventuraremos una estratexa posible pa la so entrada en campu de la tecnoloxía llingüística. Los factores qu'entren en xuegu nel accesu d'una llingua a les tecnoloxíes del llinguaxe son esencialmente les siguientes:

- *Númeru de falantes*: A mayor númberu falantes, mayor ye la rentabilidá económica que supón el desendolcu de tecnoloxía llingüística y, poro, l'interés comercial que ye a movilizar la inversión nel aria de parte d'empreses privaes. L'exemplu paradigmáticu de rentabilidá comercial ye, de xuru, l'inglés.

- *Sofitu institucional*: Nos casos nos que nun hai una rentabilidá económica visible, el sofitu y promoción de la llingua de parte de l'administración son básicos a la de garantizar el desenvolvimientu d'applicaciones de tecnoloxía llingüística. Naturalmente, el sofitu dende'l gobiernu facilita tamién el d'otres instituciones, como fundaciones culturales privaes con capacidá de mecenalgu. Nesta situación alcuéntrense llingües de pocos falantes, como'l finlandés o l'holandés.

- *Diglosia*: La entrada d'una llingua en campu tecnolóxicu depende tamién del grau nel qu'ésta sía emplegada en tolos ámbitos y rexistros d'usu, orales como escritos. Nuna situación de diglosia, una llingua supeditada y arrequexada dafechu a la comunicación oral dientro d'ámbitos familiares va tener bien difícil la creación de tecnoloxía llin-

güística: per un sitiú porque los recursos esistentes como puntu de partida (dende diccionarios y gramátiques hasta un cuerpu de testos, escritos o grabaos, abondo voluminosu) van ser mui escasos; pel otru, porque'l desendolcu d'aplicaciones va pasar prioritariamente pela llingua subordinante. Con éses, la situación d'escasez de tecnoloxía llingüística na que s'atopa l'hindi, magar el so altu númeru de falantes, seique ye debida a la so supeditación al inglés.

• *Accesu a les tecnoloxíes de la información:* Finalmente, pa qu'una llingua tenga entrada dafechu en terrén de les tecnoloxíes del llinguaxe ye mester que la so sociedá tenga tamién accesu tecnolóxicu de mou más o menos xeneralizáu. Esto nun significa la esistencia d'ordenadores personales en tolos llares de fala na llingua en cuestión (lo que supondría una visión occidente-céntrica dafechu del problema), sinón l'usu garantizáu de tecnoloxía informática nes instituciones académiques y d'investigación, mesmo que nes empreses privaes dedicaes al desendolcu de productos de base llingüística. Na mayoría de casos la imposibilidá d'accasu a les tecnoloxíes de la información vien condicionada por una seria situación de diglosia. Fai'l casu del quechua, la llingua indíxena americana más falada, con cerca de 10 millones de falantes.

Cada llingua del mundu asítiase nun puntu de coordenaes determináu en rellación a estos cuatro parámetros, lo que trai darréu un nivel determináu de capacidá de desenvolvimientu dentro de les tecnoloxíes llingüístiques. El trabayu nesta estaya puede estremase d'acordes colos tres niveles de desenvolvimientu, que corresponden, al aldu, colos plantegaos en dellos trabayos que reflexonen sobre les estratexes pal desendolcu de tecnoloxía llingüística pa llingües minoritaries (Sarasola 2000; Agirre *et al.* 2002; Diaz de Ilarraza *et al.* 2003):

• *Fundamentos:* Fase d'embalagamientu de datos léxicos y corpus testuales, entá ensin procesamientu a nengún nivel. Los productos re-

sultantes nesti estadiu tiénense de puntu de partida necesariu pal desenrueldu de tecnoloxía llingüística.

• *Ferramientes:* La bibliografía citada enantes considera esta fase como la empobinada a la igua de preseos pal procesamientu del llinguaxe natural: lematizadores, analizadores morfolóxicos y sintácticos, alliniadores de frases pa corpus multilingües, ontoloxíes o bases de conocimientu léxico-semánticu, etiquetadores semánticos, etc. Trátase de ferramientes que tanto van sirvir nel tratamientu de la información atropada na fase previa, como van valise d'ella (por exemplu, nel entrenamientu de ferramientes de base estocástica).

Amás de la fechura de ferramientes, al mio pensar, nesta fase hai qu'amestar tamién la ellaboración de preseos llingüísticos basaos nes coleccions d'información atropaes na etapa previa. Refiérome a diccionarios y gramátiques d'usu públicu. La so ellaboración en base a corpus textuales dexa espeyar de manera más afayadiza la llingua que se describe ende.

• *Aplicaciones:* Fase aplicada a la construcción d'applicaciones empobinaes a usuarios non especializaos. Por casu: correctores gramaticales, sistemas d'ayuda a la traducción con dél nivel de sofisticadura, sistemas de busca o recuperación d'información, o sistemas de diálogu. El llabor nesta etapa mándase necesariamente de les ferramientes desendolcaes na fase anterior.

Pa les llingües con un númberu curtiu de falantes, en clara situación de diglosia y ensin sofitu institucional, la entrada na tecnoloxía de la información va ser, si non imposible, costosa de manera. Ésta ye, llamentablemente, la situación de la inmensa mayoría de les 6800 llingües qu'entá anguaño se falen nel nuesu mundu. Puestos no mejor, el so contactu cola tecnoloxía llingüística va ser pentemedies de la Llingüística de Corpus, si ye que daqué grupu investigador s'alcuerda d'entamar el proyeetu de recopilación de materiales léxicos ya igua de cor-

pus (na mayoría de casos, ensin nengún tipu de procesamientu nin anotación) enantes de que la llingua en cuestión desapaeza<sup>16</sup>.

Nun asitiamientu intermediu afáyense, per un sitiu, les llingües pequeñes con sofitu institucional, con usu xeneralizáu a tolos niveles y una capacidá tecnolóxica prestosa, tal como holandés y finlandés. Y per otru, les grandes llingües con un grau d'accusu tecnolóxicu relativamente baxu que torgara apocayá la so entrada en campu de les tecnoloxíes llingüístiques, pero que por razones comerciales o estratéxiques tiren agora a xenerar ciertu interés. Fai'l casu del farsi o'l tagalog. Nesti segundu nivel, el desenrueldu de tecnoloxía llingüística resal del nivel primariu d'atropamiento de corpus ensin etiquetar, por entrar yá na etapa de desenvolvimientu de ferramentes básiques pal procesamiento de corpus, o la igua de preseos de base llingüística (diccionarios y gramátiques) que pudieren mandase de corpus testuales básicos y bases de datos léxicos embalagaos na etapa previa.

D'últimes, atopamos les llingües que presenten un gran número de falantes y que gocen de lo que calificara como un contestu de fácil accesu tecnolóxicu, como l'inglés, el chinu o l'español, daveres la ínfima minoría. Pa estos casos, el nivel de desenrueldu de les tecnoloxíes del llinguaxe ta al llegar, si entá nun lo fexo, a la creación d'applicaciones empobinaes a usuarios non expertos como los mentaos na tercer etapa del desendolcu de la tecnoloxía llingüística.

Naturalmente, l'estremar en tres grandes bloques les llingües del mundu d'acordies col so grau de participación nes tecnoloxíes del llinguaxe ye la simplificación d'una situación enllena de matices. Un exemplu d'ello failu'l casu de les llingües minoritaries más cercanes al

<sup>16</sup> A esto dedícase por casu *SIL International* (ente otros proyectos; véase: <http://www.sil.org>) o programes de mecenalgu como'l *Documentation Programme* del *The Hans Rausing Endangered Languages Project* ([http://www.hrelp.org/doc\\_home.htm](http://www.hrelp.org/doc_home.htm)).

asturianu: el catalán, el euskera y el gallegu. Trátase de tres llingües con un número reducíu de falantes (unos 6 millones, 500.000 y 3,5 millones respectivamente de falantes como primer llingua) y so la influencia de daqué nivel de diglosia respectu al español (mayor o menor, según el casu), amás d'enguedeyos derivaos d'una fuerte variación dialectal nel euskera y, bien qu'en menor grau, en gallegu. Sicasí, el desenrueldu de tecnoloxía llingüística nestes tres llingües ye de sorrayase: nos tres casos se degolara la fase inicial de creación de los fundamentos y tiense entrao yá na segunda etapa, cola consiguiente xeneración de productos de base llingüística (como'l *Diccionari de freqüències del catalán*, o bien el *XX Mendeko Euskararen Corpus Estadistikoa* pal euskera, consultable per internet<sup>17</sup>), y na construcción de ferramientes de procesamientu del llinguaxe, indispensables pa la entrada nel nivel de les aplicaciones d'usuariu non especializáu. Amás, tanto en casu del euskera como nel del catalán, comercializarónse yá delles aplicaciones finales correspondientes al tercer nivel de desendolcu (véase Diaz de Ilarrazu *et al.* [2003] pal euskera, y Llisterri [2000] pal catalán). Nun dispongo d'información nesti sen no que fai al gallegu.

Cuento qu'esta situación respuende principalmente a dos factores. Per un sitiu, el compromisu cola llingua autóctona y l'enfotu de trabayar pola so normalización dende les esferes académiques y d'investigación. Otra manera, la voluntá de les administraciones nacionales catalana, vasca y gallega d'afalar el trabayu nel área, por cuantes se camienta que namá d'esta miente puede asegurase la competitividá y validez de la llingua nel nuevu marcu de la sociedá de la información<sup>18</sup>. A es-

---

<sup>17</sup> Pa ún y otru trabayu, ver les referencies nes correspondientes secciones.

<sup>18</sup> A éstes, ye relevante la propuesta d'acción estratéxica nes aries d'Industria de la Llingua ya Inxeniería Llingüística de parte del gobiernu vascu (véase la nota 11 pa la referencia completa), o bien la creación del *Centre de Referència en Enginyeria Lingüística* (CREL) pola Generalitat de Catalunya, y el *Centro Ramón Piñeiro para a Investigación en Humanidades* de parte de la

tos dos elementos hai que-yos amestar un terceru: el nivel de desenvolvimientu de les tecnoloxíes llingüístiques de qu'esfruten estes llingües, magar la so condición, namá ye posible pa les llingües minoritaries del mundu occidental, cuidao que tienen un mínimu d'estabilidá económica garantizáu, accesu a les tecnoloxíes de la información, y tán quites d'un contestu de marxinación social como ye'l casu de les llingües indíxenes en tolos países ensin escepción del continente americanu. L'estáu de desenvolvimientu de la tecnoloxía llingüística pa les neses tres llingües nun ye óptimu a comparalu cola situación del inglés o l'español, claramente valíes pola inversión del sector priváu. Si casí, ye bien superior al estáu de la mayoría de llingües del mundu que s'alcuentren nuna situación comparable.

Na redolada europea, l'aplicación de les tecnoloxíes del llinguaxe nel ámbitu de les llingües minoritaries autóctones ye anguaño un tema de creciente interés. Amuesa d'ello ye la pasu ente pasu más bayu-rosa bibliografía sobre'l tema, la creación apocayá de SALTMIL<sup>19</sup>, un grupu d'interés dedicáu a ello (Nadeu *et al.* 2001), o bien los seminarios temáticos entamaos en marcu de congresos d'ampliu algame que se centraron nesta cuestión<sup>20</sup>. La premisa de partida ye l'aceptación de la llamada Sociedá de la Información como'l nuevu contestu de referencia pal truecu de conocimientu y, darréu, l'avance tecnolóxicu d'occidente, neto que'l papel indispensable de la Inxenie-

---

Xunta de Galicia, dos instituciones aplicaes a la investigación y desenrueldu de tecnoloxía llingüística pa les llingües autóctones.

<sup>19</sup> SALTMIL ye'l Special Interest Group on Speech and Language Technology for Minority Languages, creau dientro de la International Speech Communication Association (ISCA) (<http://isl.nftex.uni-lj.si/SALTMIL/>).

<sup>20</sup> Por exemplu los workshops: «Language Resources for European Minority Languages» (Granada, LREC 1998), «Developing language resources for minority languages: re-useability and strategic priorities» (Atenas, LREC 2000), y «Natural Language Processing Of Minority Languages And Small Languages» (Batz-sur-Mer, France, TALN 2003).

ría Llingüística pal so progresu. Con éstes, la supervivencia de les llingües minoritaries occidentales pasa necesariamente pel desendolcu de tecnoloxía llingüística (Sarasola 2000, Nadeu *et al.* 2001).

Tiene visto que la complexón vital de cada llingua determina a grandes rasgos la so posibilidá de desenvolverse dentro de les tecnoloxías del llinguaxe. Sicasí, tien que se considerar tamién que la entrada d'una llingua minorizada nesti campu supón, arriendes de la posibilidá d'actualización y ameyoramientu de la so presea de descripción llingüística (diccionarios y gramátiques), la so entrada en contestu de la sociedá de la información, que ye lo de calificala como llingua moderna que ye quien a competir a nivel tecnolóxicu, abandonando la so imaxe arcaica, ensin cohesión interna nin capacidá pa la comunicación a tolos niveles. En definitiva, contribúi al frayamientu de la situación de diglosia na que s'afaya, lo que, otramiente, retroalimenta positivamente la voluntá de sofitu institucional y aguiya un progresivu interés comercial<sup>21</sup>. El desenrueldu de les tecnoloxíes del llinguaxe ye, poro, daqué que va tresallá d'un cenciellu exerciciu académicu. Ye un asuntu d'interés públicu y de sonadía social, nel de xugase l'afitamientu y supervivencia de les llingües minoritaries como llingües vives nuna cultura global de base tecnolóxica.

### *Puntu de partida pal asturianu*

La entrada del asturianu na tecnoloxía llingüística vuélvese, con éses, un proyeetu de primer prioridá, oldeable a la so incorporación

<sup>21</sup> Nesti sen abúltame ilustrativa una anécdota tocántenes a la entrada del catalán na industria televisiva y cinematográfica: de la qu'en 1983 estrenóse TV3, la canal de televisión de Cataluña, sentir falar en catalán a *JR, Sue Ellen* y los sos compañeros de la serie *Dallas* (yá conocíos del público al traviés de TVE) yera motivu de risión. En 2002 sicasí, la indignación popular y el so consiguiente boicot comercial, movíos pola refuga de la *Warner Bros* a doblar *Harry Potter* al catalán, fexo camudar la opinión de la multinacional norteamericana.

nos medios de comunicación de masas. La coyuntura política actual paez abrir una posibilidá al proyectu: de parte de los organismos responsables del aria na Comunidá Europea, esiste una reconocencia de la situación que se-yos plantega anguaño a les llingües minoritaries autóctones nesti nuevu marcu d'información.

Amás, la llingua asturiana esfruta yá d'una mínima base de partida. De mano, dispón de l'Academia de la Llingua Asturiana, una institución con autoridá llingüística, fundamental p'articular proyectos d'investigación. De segundes, tien entrao nuna fase de normalización gracias a la publicación de la *Gramática de la Llingua Asturiana* (3<sup>a</sup> ed., 2001) y del *Diccionariu de la Llingua Asturiana* (2000) a cuenta d'esta mesma academia. D'últimes, dispón d'un decente cuerpu de testos escritos, sobremanera nos rexistros lliterarios y periodísticos, de los qu'una parte foi yá dixitalizada por aciu del *Proyectu Caveda y Nava*<sup>22</sup>, mesmo qu'un archivu oral –bien que necesariamente ampliable – atropáu en marcu del proyectu *Archivu Oral de la Llingua Asturiana*<sup>23</sup>.

Sicasí, nun tien que se faer de menos la situación enxeble na qu'anguaño s'atopa la llingua. El so estatus de non oficialidá traduzse na práctica nuna menor capacidá a la de beneficiase de sofitu económico. Amás, esiste la constatación d'una perda xeneracional de falantes, neto qu'un amenorgamientu del usu de la llingua mientres la última década en rexistros indagora eslusivos d'ella (Llera Ramo 2003). Estos dos factores suponen un acutamiento considerable en términos de recursos, y, poro, la propuesta pa un proyectu que dea entrada al asturianu nes tecnoloxíes de la información nun puede entamar per plantegamientos maximalistes. Otramiente, hai que s'atener a la experiencia sacada de proyeutos con otros llingües en situación asemeyada, como l'euskera.

<sup>22</sup> <http://www.cavedaynava.org/> (URL vixente el 29 d'ochobre del 2003).

<sup>23</sup> <http://www.asturias.org/asturianu/archoral/> (URL vixente el 20 d'ochobre del 2003).

Nestos trabayos previos viose cómo nun ye d'empezar a desenroldar ferramientes nin aplicaciones si nun esisten los fundamentos; esto ye, compilaciones léxiques y corpus testuales (Sarasola 2000, ente otros).

Aplicáu al casu particular de la llingua asturiana, esto pica a la creación d'un corpus testual –escritu y/o oral– como prioridá inicial pa la so entrada en campu de les tecnoloxíes llingüístiques. Un corpus testual, en xunto col atropamiento de datos léxicos (que puede tener el diccionariu de l'Academia de la Llingua Asturiana como puntu de partida), garantiza l'encontu pal desenvolvimientu posterior de ferramientes y aplicaciones. De la mesma, permite a mediu plazu la fechuura de productos de base llingüística. Por casu, los estudios llingüísticos realizaos sobre los datos del corpus pueden derivar n'aplicaciones como l'ameyoramientu de la gramática y el diccionariu mentaos enantes, la creación d'un diccionariu descriptivu a cuenta del léxicu más frecuente de la llingua, o bien la planificación de delles actuaciones sociollingüísticas. Otramiente, ye posible pensar n'aplicaciones que valgan al aprendizaxe del asturianu: cuándo ameyorando los materiales esistentes a partir de datos procedentes del corpus, cuándo desenrol dando ferramientes pal aprendizaxe valíu por ordenador –por exemplu, exercicios que requieran reproducir les categoríes morfosintáctiques asignaes en determinaes frases del corpus, xenerar formes flexionaes a partir d'un lema y la información a él venceyada, completar perífrasis, frases feches y allugamientos, etc.

Nenguna d'estes aplicaciones quier un nivel de tecnoloxía sofistícau por demás. Un corpus testual lematizáu y etiquetáu morfosintácticamente, en xunto a una ferramienta d'esplotación de los datos, fai abondo. Asitiámonos namá nel segundu nivel de desenrueldu tecnolóxicu presentáu na sección, y, sicasí, les posibilidaes de les aplicaciones que se planteguen son yá notables.

Repárese, con eso, que se trata d'aplicaciones que presuponen un fragmentu particular de lo que ye la llingua asturiana: l'asturianu ac-

tual. La entrada del asturianu nes tecnoloxíes del llinguaxe plantégase (arriendes de pol inherente interés a nivel científicu) en ficies a derrompe-y la entrada al nuevu marcu de comunicación y llexitimalu en cuantes que llingua válida en tolos niveles d'usu. Teniendo de cuenta la so presente situación, un corpus de llingua actual abulta un beneficiu inicial más xeneralizable a distintes aries de la llingüística aplicada que non, por casu, un corpus históricu.

Otra manera, ye relevante de manera qu'esti primer corpus potential recueya la máxima variedá dialectal, diastrática y diafásica del asturianu, por mor del estadiu d'asimilación de la norma estándar nel que s'alcuentra anguaño, a pocos años de l'apaición d'un diccionariu y una gramática asoleyaoa pola institución con autoridá normativa sobre la llingua. Cuandoquier, interesa qu'esti primer corpus de la llingua asturiana sia un corpus de referencia de llingua xeneral.

Un corpus pal asturianu tien que se concebir tamién con cuenta de que, a la llarga, se mande como base pal desenrueldu d'applicaciones tecnolóxiques más sofisticadas: sistemas de diálogo, sistemas de summarización, de recuperación o cata d'información, etc. D'ésta, son mester dos elementos. D'uno, que les ferramientes básiques de procesamientu, como lematizadores, analizadores morfosintácticos o sintácticos, se fixeren de mou que puedan ser bonos de reutilizar, a faese necesario, de parte de caúna d'estes aplicaciones finales (nuevamente Sarasola 2000 y Diaz de Ilarrazo *et al.* 2003 defenden esti criteriu, resultanza de la so experiencia col euskera). De segundes, que la igua d'estes aplicaciones finales surda d'un análisis de la realidá del asturianu y la so adecuación a ella.

Esto ye: chando cuenta de la condición de llingua minoritaria del asturianu, nun ye dable plantegar el desendolcu de toles aplicaciones nes que se trabaya anguaño pa llingües como l'inglés. Pela cueta, hai que concentrar l'esfuerzu naquelles xeres pa les qu'esista garantía real d'usu –y, con ella, un mínimo de rentabilidá comercial–. Considerando les condiciones de la llingua asturiana, abulta claro que, por exem-

plu, un sistema de sumarización de documentos va a ser d'un interés y nivel d'aplicación en forma inferior al que puede tener un sistema de traducción automática ente l'asturianu y l'español. Poro, el proyectu de construcción d'un corpus de llingua asturiana tien de considerar, na mio opinión, la posibilidá de que siquier una parte sea billingüe asturiano-español, preferentemente paralelu.

Finalmente, el proyectu tamién tien de char cuenta del estáu nel que s'afaya el fragmentu de llingua que mire a representase. Por exemplu, va tener que tomar decisiones sobre'l tratamiento de la variación dialectal; por exemplu, la qu'afecta a la flexón de determinaes categoríes, pola repercusión qu'esto puede tener en procesu de lematización y asignación de etiquetes morfosintáctiques. O, otra manera, va tener qu'establecer una política específica pa la variación ortográfica, a la de tratar con testos d'una etapa pre-normativa. Nesti sen, puede sirvir de referencia la esperiencia d'otros proyectos desendolcaos pa llingües cercanes al asturianu, como les citaes.

#### BIBLIOGRAFÍA CITADA

ABAITUÁ (1999)= JOSEBA ABAITUÁ, «Quince años de traducción automática en España», *Perspectives: Studies in Translatology*, 7-2 (1999), páxs 221-230. Versión catalana: «Quinze anys de traducció automàtica a l'Estat espanyol», Digit.HUM. Universitat Oberta de Catalunya. <http://www.uoc.es/bumfil/digitum/digitum2/catala/didactica/index.html>.

AGUIRRE *et al.* 1998 = E. AGUIRRE, I. ALDEZABAL, I. ALEGRIA, O. ANSA , X. ARREGI, J. ARRIOLA, X. ARTOLA, A. DÍAZ DE ILARRA, N. EZEIZA, K. GOJENOLA , A. MARITXALAR, M. MARITXALAR, M. ORONoz, K. SARASOLA, A. SOROA, R. URIZAR, M. URKIA (1998), «A framework for the automatic processing of Basque», *Workshop on Lexical Resources for Minority Languages*. Granada (Language Resources and Evaluation Conference), 1998.

AGUIRRE *et al.* 2002 = E. AGUIRRE, I. ALDEZABAL, I. ALEGRIA, X. ARREGI, J. M. ARRIOLA, X. ARTOLA, A. DÍAZ DE ILARRAZA, N. EZEIZA, K. GOJENOLA, K. SARASOLA, A. SOROA, «Towards the definition of a basic toolkit for HLT», LREC 2002. *Third International Conference On Language Resources And Evaluation*. Las Palmas (Language Resources and Evaluation Conference), 2001.

AGUIRRE MORENO *et al.* 2001 = J. L. AGUIRRE MORENO, N. ANDIÓN, X. GÓMEZ GUINOVART, «Aspectos ortográficos, léxicos y morfosintácticos del etiquetado lingüístico de un corpus de informática en lengua gallega», *Procesamiento del Lenguaje Natural*, 27 (2002), páxs. 13-19.

AGUIRRE MORENO *et al.* (2002) = J. L. AGUIRRE MORENO, A. ÁLVAREZ LUGRÍS, X. GÓMEZ GUINOVART, «Etiquetario morfosintáctico del SLI para corpus de lengua gallega: aplicación al corpus paralelo TECTRA», *Procesamiento del Lenguaje Natural*, 28 (2002), páxs. 23-34.

AGUIRRE MORENO *et al.* (2003) = J. L. AGUIRRE MORENO, A. ÁLVAREZ LUGRÍS, I. BRAGADO TRIGO, L. CASTRO PENA, X. GÓMEZ GUINOVART, A. LÓPEZ LÓPEZ, J. R. PICHEL CAMPOS, E. SACAU FONTENLA, L. SANTOS SUÁREZ «Alinhamento e etiquetagem de corpora paralelos no CLUVI (Corpus Linguístico da Universidade de Vigo)», J. J. ALMEIDA (ed.), *Workshop CP3A 2003, Corpora Paralelos: Aplicações e Algoritmos Associados*. Braga (Universidade do Minho), 2003.

ASTON y BURNARD 1998 = GUY ASTON, LOU BURNARD, *The BNC Handbook. Exploring the British National Corpus with SARA*, Edinburgh (Edinburgh University Press), 1998.

BACH *et al.* 1997 = CARME BACH, ROSER SAURÍ, JORDI VIVALDI, M. TERESA CABRÉ, *El Corpus de l'IUJA: descripció*, Barcelona (Institut Universitari de Lingüística Aplicada, Serie Informes), 1997.

BIBER 1993 = D. BIBER, «Representativeness in corpus design», *Literary and Linguistic Computing*, vol. 8 (4), 1993, páxs. 243-257.

BOIX 1996 = E. BOIX, «Els materials de llengua oral dels corpora de català contemporani de la UB (CUB)», LL. PAYRATÓ, E. BOIX, M. R. LLORET, M. LORENTE (eds.) *Corpus, Corpora. Actes del 1er i 2on Col.loquis Lingüístics de la Universitat de Barcelona (CLUB-1, CLUB-2)*, Barcelona (Promociones y Publicaciones Universitarias S A), 1996, páxs. 93-114.

CARLSON *et al.* 2003 = LINN CARLSON, DANIEL MARCU, MARY ELLEN OKUROWSKI, «Building a discourse-tagged corpus in the framework of rhetorical structured theory», JAN VAN KUPPEVELT, ROONIE SMITH (eds.) *Current Directions in Discourse and Dialogue*, Kluwer (en prensa), 2003.

DÍAZ DE ILARAZA *et al.* 2003 = A. DÍAZ DE ILARAZA, A. GURRUTXAGA, I. HERNÁEZ, N. LÓPEZ DE GEREÑU, K. SARASOLA, «HIZKING2I: Integrating language engineering resources and tools into systems with linguistic capabilities», *Workshop on NLP of Minority Languages and Small Languages*, Batz-sur-Mer (Traitement Automatique des Langues Naturelles), 2003.

FLIGELSTONE 1992 = STEVE FLIGELSTONE, «Developing a scheme for annotating

text to show anaphoric relations», Gerhard Leitner (ed.), *New Directions in English Language Corpora: Methodology, Results, Software Developments*, Berlin (Mouton de Gruyter), 1992.

FRANCIS y KUCERA 1964 = W. N. FRANCIS, H. KUCERA, *Manual of Information to Accompany a Standard Sample of Present-Day Edited American English, for use with Digital Computers*. Providence RI (Brown University, Department of Linguistics. Edición revisada y ampliada), 1979. <http://www.hit.uib.no/icame/brown/bcm.html> (última actualización: 11/9/1997).

GARCÍA-MATEO *et al.* 1998 = CARMEN GARCÍA-MATEO, MANUEL GONZÁLEZ-GONZÁLEZ, «An overview of the existing language resources for Galician». *Workshop on Language Resources for European Minority Languages*, Granada (Language Resources and Evaluation Conferences), 1998.

LLERA RAMO 2003 = FRANCISCO JOSÉ LLERA RAMO, «Llumes y solombres na evolución sociollingüística del asturianu», *XXII Xornadas Internacionais d'Estudio de la Llingua Asturiana*, Academia de la Llingua Asturiana, Uviéu, 27-29 ochobre 2003.

LLISTERRI 1997 = JOAQUIM LLISTERRI, «Transcripción, etiquetado y codificación de corpus orales», *Seminario de Industrias de la Lengua*, Soria (Fundación Duques de Soria), 1997.

LLISTERRI 1999 = JOAQUIM LLISTERRI, «Corpus orals per a la fonètica i les tecnologies de la parla», *Actes del I Congrés de Fonètica Experimental*. Tarragona, 22, 23 i 24 de febrer de 1999, Tarragona (Universitat Rovira i Virgili - Universitat de Barcelona), páxs. 27-38.

LLISTERRI 2000 = JOAQUIM LLISTERRI, «O catalán nas industrias da lingua», *Lingua e cultura catalanas*, Vigo (Universidade de Vigo, Cursos de Extensión Universitaria), 2000.

MCENERY y WILSON 2001 = TONY MCENERY, ANDREW WILSON, *Corpus Linguistics*, Edinburgh (Edinburgh University Press), 2001.

MANN y THOMPSON 1988 = WILLIAM MANN, SANDRA THOMPSON, «Rethorical Structure Theory. Toward a functional theory of text organization», *Text*, 8/3 (1988), páxs. 243-281.

NADEU *et al.* 2001 = CLIMENT NADEU, DONNCHA ÓCRÓINÍN, BOJAN PETEK, KEPÀ SARASOLA, BRIONY WILLIAMS, ISCA SALTMIL SIG, «Speech and language technology for minority languages», <http://gps-tsc.upc.es/veu/research/pubs/download/Nado1b.pdf> (28/10/2003).

PINO *et al.* 1999 = M. PINO, M. SÁNCHEZ SÁNCHEZ, «El subcorpus del Banco de Datos CREA-CORDE (RAE): procedimientos de transcripción y codificación», *Oralia*, 2 (1999), páxs. 83-138.

PUSTEJOVSKY 2003 = JAMES PUSTEJOVSKY, PATRICK HANKS, ROSER SAURÍ, ANDREW SEE, ROBERT GAIZAUSKAS, ANDREA SETZER, BETH SUNDHEIM, LISA FERRO, «The TIMEBANK Corpus», *Proceedings of the Corpus Linguistics 2003 Conference*, Lancaster (Lancaster University Press), 2003.

QUIRK 1992 = R. QUIRK, «On corpus principles and design», JAN SVARTVIK (ed.) (1992), *The London Corpus of Spoken English: Description and Research*, Lund (Lund University Press), 1992, páxs. 457- 469.

QUIRK 1985 = R. QUIRK, S. GREENBAUM, G. LEECH & J. SVARTVIK, *A comprehensive Grammar of the English Language*, Harlow (Longman), 1985.

RAFEL 1992-93 = J. RAFEL, «El “Diccionari del català contemporani”: Treballs realitzats i previsions de futur», *Llengua i Literatura*, 5 (1992-93), páxs. 733-737.

RAFEL 1996 = J. RAFEL, «El Diccionari del català contemporani i el Corpus textual informatitzat de la llengua catalana», LL. PAYRATÓ, E. BOIX, M. R. LLORET, M. LORENTE (eds.) *Corpus, Corpora. Actes del 1er i 2on Col·loquis Lingüístics de la Universitat de Barcelona (CLUB-1, CLUB-2)*, Barcelona (Promociones y Publicaciones Universitarias SA), 1996, páxs. 71-92.

RAFEL I FONTANALS 1996-98 = J. RAFEL I FONTANALS (dir.), *Diccionari de freqüències*, 3 vols., Barcelona (Institut d'Estudis Catalans), CD-ROM de la obra completa, 1996-98.

RODRÍGUEZ 2003 = CARLOS RODRÍGUEZ, «Applying information extraction techniques to metalinguistic discourse», *Topics in Computational Linguistics and Intelligent Text Processing*, Norwell (Springer Verlag, Lecture Notes in Computer Science), 2003.

SÁNCHEZ *et al.* 1995 = A. SÁNCHEZ, R. SARMIENTO, P. CANTOS, J. SIMÓN, *Cumbré. Corpus lingüístico del español contemporáneo. Fundamentos, metodología y aplicaciones*, Madrid (SGEL), 1995.

SÁNCHEZ-LEÓN *et al.* 1999 = F. SÁNCHEZ-LEÓN, J. PORTA, J. L. SANCHO, A. NIETO, A. BALLESTER, A. FERNÁNDEZ, J. GÓMEZ, L. GÓMEZ, E. RAIGAL, R. RUIZ, «La anotación de los corpus CREA y CORDE», *Actas del XV Congreso de la Sociedad Española para el Procesamiento del Lenguaje Natural*, MADRID (Sociedad Española para el Procesamiento del Lenguaje Natural), 1999.

SARASOLA 2000 = K. SARASOLA, «Strategic priorities for the development of language technology in minority languages», *Workshop on Developing Language Resources for Minority Languages: Re-useability and Strategic Priorities*, Athens (Language Resources and Evaluation Conference), 2000.

SINCLAIR 1987 = J. SINCLAIR (ed.), *Looking Up, An Account of the COBUILD Project*, London (Collins), 1987.

SINCLAIR 1986 = J. SINCLAIR, *Preliminary recommendations on corpus typology*. EA-GLES Document TCWG-CTYP/P. <http://www.ilc.cnr.it/EAGLES96/corpustyp/corpustyp.html> (21/9/03), 1986.

SOLER I BOU 1998 = J. SOLER I BOU, «Los corpus textuales en lengua catalana», *Curso de industrias de la lengua. Proyectos actuales en procesamiento de lenguaje natural*, Soria (Fundación Duques de Soria), 1998.

SOMERS 2001 = HAROLD SOMERS, «Where do we stand?», Panel Session, *MT Summit VIII*, Santiago de Compostela, 18-22, september 2001.

SVARTVIK (ed.) 1990 = JAN SVARTVIK, *The London Corpus of Spoken English: Description and Research*. Lund Studies in English 82, Lund (Lund University Press; Lund Studies in English 82), 1990.

TOGNINI-BONELI 1996 = E. TOGNINI-BONELI, *Corpus Theory and Practice*, Birmingham (Tuscan Word Center), 1996.

TORRUELLA Y LLISTERRI 1999 = JOAN TORRUELA, JOAQUIM LLISTERRI, «Diseño de corpus textuales y orales». J. M. BLECUA, G. CLAVERÍA, C. SÁNCHEZ, J. TORRUELLA (eds.) *Filología e informática. Nuevas tecnologías en los estudios filológicos*, Barcelona (Universitat Autònoma de Barcelona, Editorial Milenio), 1999, páxs. 45-77.

VIAPLANA 2000 = J. VIAPLANA, «Corpus oral de variació», *Jornades del Centre de Referència en Enginyeria Lingüística (CREL)*, Institut d'Estudis Catalans, Barcelona, 4 i 5 d'abril de 2001.