

# Syntactic annotation in the Reference Corpus for the Processing of Basque (EPEC): Theoretical and practical issues\*

IZASKUN ALDEZABAL, MARIA JESUS ARANZABE,  
JOSE MARI ARRIOLA, and ARANTZA DIAZ DE ILARRAZA

## 1 *Abstract*

2 *In this paper, we will describe some theoretical and practical issues raised*  
3 *during the construction of the Basque Dependency Treebank (BDT): the syn-*  
4 *tactic annotation of EPEC (Reference Corpus for the Processing of Basque).*  
5 *EPEC is a 300,000 word corpus of standard written Basque whose purpose*  
6 *is to be a training corpus for the development and improvement of several*  
7 *NLP (Natural Language Processing) tools for Basque. BDT will be the first*  
8 *corpus for the Basque language tagged at syntactic level. We will also present*  
9 *the dependency-based annotation hierarchy that we have established for the*  
10 *syntactic tagging. Decisions made during design of the annotation hier-*  
11 *archy are based on the description of Basque grammar made by Euskaltzaindia*  
12 *(Academy for the Basque Language). When describing dependency relations,*  
13 *we consider lexical units as syntactic heads. This will open up a way for us*  
14 *to work with semantics.*

15 *Keywords: PLEASE ADD!*

## 16 **1. Introduction**

17 A treebank is a text corpus in which each sentence has been annotated with  
18 its syntactic structure. The construction of a treebank is a multidisciplinary  
19 task that, although expensive, is indispensable for the development of real  
20 applications in the field of NLP. At a purely linguistic level, a Treebank  
21 is an essential database for the study of a language given that it provides  
22 analyzed/annotated examples of real language. Besides, the linguistic study  
23 produces an improvement in the quality of several applications, such as Part-  
24 Of-Speech (POS) taggers and parsers (Collins 1997, 2000; Charniak 2000),  
25 because it provides common training and testing material allowing different  
26 algorithms to be compared and improved.

27 Over recent years, treebank corpora such as the Penn Treebank (Marcus  
28 et al. 1993) and the Prague Dependency Treebank (Böhmová et al. 2003)  
29 have become a crucial resource for building and evaluating natural language  
30 processing tools and applications. Abeillé (2003) explained the work for  
31 Czech, German, French, Japanese, Polish, Spanish and Turkish, to name just  
32 a few. Kakkonen (2005) presents the state of the art of dependency-based  
33 treebanks.

34 The Basque Dependency Treebank (BDT) is currently the Reference Cor-  
35 pus for the Processing of Basque (EPEC), annotated at syntactic level. EPEC  
36 is a 300,000 word corpus of standard written texts which is intended to be a  
37 training corpus for the development and improvement of several NLP tools  
38 (Bengoetxea and Gojenola 2006).

39 In this paper, we describe the theoretical and practical issues raised dur-  
40 ing construction of the BDT, following the Dependency Grammar theory  
41 (Tesnière 1959). This is the first formalization of the syntactic tagging of  
42 Basque that follows the Dependency Model; it should be noted that we have  
43 based our work on the syntactic description of Basque grammar made in  
44 Euskaltzaindia (1991 [1985], 1987, 1990, 1994, 1999). Using dependency  
45 relations we have formalized the main syntactic structures described in this  
46 grammar. We have also made our own decisions during the design of the  
47 annotation hierarchy. This is very important, as Sampson (2003: 23–41) says  
48 with respect to the SUSSANE corpus: “In the work of my group I have cho-  
49 sen the opposite priority: we treat the detail, accuracy and explicitness of  
50 annotation as more important than the quantity of material annotated, with  
51 the inevitable consequence that our treebanks have to be very small”.

52 Dependency theory is one of the most widely used methods of conceptu-  
53 alizing the linguistic structure of sentences. In grammars constructed using  
54 the dependency theory (Hudson 1990; Mel’cuk 1988), syntax is handled in  
55 terms of grammatical relations between pairs of individual words, such as  
56 the relation between the subject and the predicate or between a modifier and  
57 a common noun. Grammatical relations are seen as subtypes of a general,  
58 asymmetrical dependency relation: one of the words (the head) determines the  
59 syntactic and semantic features of the combination. The syntactic structure of  
60 a sentence as a whole is built up from the dependency relations between indi-  
61 vidual pairs of words. In general terms, we take as syntactic heads the lexical  
62 elements that are involved in the dependency relations. The decision to take  
63 lexical units as a basis for any dependency relation is also motivated by the  
64 previously developed syntactic approach: the dependency-oriented surface  
65 syntax (Järvinen and Tapanainen 1997). This will allow us to extend our work  
66 to semantics in the near future. A similar approach is followed by Lin (1995).

67 Taking into account the literature on tagging corpora for different lan-  
68 guages and the fact that Basque syntax has been mainly developed within the  
69 generative framework of Goenaga (1991), Eguzkitza (1993), Laka (1993), Ar-

70 tiagoitia (2002), Trask (2003) and Zabala (2003), we decided to focus on the  
 71 specification and development of the annotation scheme to build the Basque  
 72 Treebank (BDT) without attempting to justify or elaborate in any depth on  
 73 any theory. Our approach is intended to provide consistent annotation to fa-  
 74 cilitate automatic exploration of linguistic data. Indeed, when designing the  
 75 annotation schema we do not have any linguistic theory in mind, apart from  
 76 Dependency Grammar Theory, so that, depending on the phenomena we have  
 77 to deal with, we determine the most accurate relation tag; for instance, in the  
 78 case of elided elements such as *pro*<sup>1</sup> we adopt the generative approach in order  
 79 to make a deeper syntactic analysis (the main peculiarities of our annotation  
 80 practice are presented in Section 4).

81 The remainder of this paper is organized as follows: Section 2 presents  
 82 the general features of the EPEC Corpus; in Section 3 we will set out the  
 83 main decisions taken in developing the syntactic annotation scheme; Section  
 84 4 describes, using examples, the annotation decisions made; and finally, some  
 85 conclusions and future work are outlined in Section 5.

## 86 2. Description of the corpus

87 The EPEC Corpus of written Basque is a 300,000 word collection of written  
 88 standard Basque. It is intended to be a reference corpus for the development  
 89 and improvement of several NLP tools for Basque. A small part of this collec-  
 90 tion has been obtained from the EEBS project (<http://www.euskaracorpora.net>), and the rest from Euskaldunon Egunkaria (<http://www.egunero.info>),  
 91 the only daily newspaper written entirely in standard Basque, published in  
 92 the second half of 1999 and in 2000. The articles were chosen to cover an  
 93 assorted range of topics (economics, culture, international, local, opinion,  
 94 politics, sports, entertainment . . .). This corpus is being used for Natural  
 95 Language Processing and although it is small, it is a strategic resource for a  
 96 minority language like Basque.  
 97

98 The corpus has been linguistically annotated at different levels: it was first  
 99 morphologically analyzed by means of MORFEUS (Alegria et al. 1997) and  
 100 then manually disambiguated (Aldezabal et al. 2007a). In the manual tagging,  
 101 each word form of the whole corpus was assigned its corresponding analysis  
 102 at the segmentation level: part-of-speech, number, definiteness and declen-  
 103 sion case. After the morphological disambiguation, other modules within  
 104 the IXATI chunker (Alegria et al. 2006; Aduriz et al. 2006), such as complex  
 105 postpositions, name entities, multiword lexical units and morphosyntax, were  
 106 applied. The manual dependency-based syntactic annotation started at this  
 107 stage. Thus, nowadays we have a Treebank for Basque of 300,000 words com-  
 108 pletely and correctly analyzed at dependency level (Aldezabal et al. 2007b;  
 109 Aranzabe 2008).



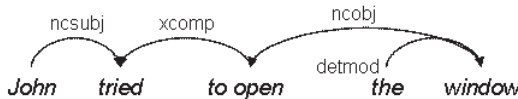
- 148 This method has three outstanding properties:
- 149 – It is based on linear word order; i.e. the order of syntactic components
  - 150 reflects the order in which they appear in the sentence.
  - 151 – Hierarchical information is made explicit.
  - 152 – The information function must be inferred.

153 *3.2. Dependency-based formalism*

154 Unlike the constituency-based approach, dependency-based formalism  
 155 (Järvinen and Tapanainen 1997) describes the relations between the com-  
 156 ponents.

157 This tagging formalism has been used for German (NEGRA) (Brants et al.  
 158 2003) and Czech (PDT) corpora<sup>3</sup>, among others.

159 In this formalism, the representation of the sentence “John tried to open  
 160 the window” above would be as follows:



- 162 The features of this method are:
- 163 – The relevance of word order is minimized.
  - 164 – It is strongly based on hierarchical relations.
  - 165 – The functional information is extremely important.

166 *3.3. Constituency-based vs. dependency-based formalism*

167 There is still an ongoing debate as to whether a constituency-based or a  
 168 dependency-based formalism should be employed in completing the tree-  
 169 bank. Some researchers have taken the middle-ground between these two op-  
 170 tions, as in Montemagni et al. (2003), who employ the dependency-based ap-  
 171 proach only to combine the basic components of the sentence (noun phrases,  
 172 prepositional phrases and the verb), without reaching the word-level for de-  
 173 pendency purposes.

174 The formalisms described above may be generally suitable, but the success  
 175 and influence they may exert on applications depends to a great extent on the  
 176 language under consideration. Based on a number of tests, set out in Skut  
 177 et al. (1997), Tapanainen and Järvinen (1998) and Oflazer et al. (1999), to  
 178 deal with the free word order displayed by Basque syntax, we have decided  
 179 to follow the dependency-based procedure. The following issues also had a  
 180 critical influence on our decision:

- 181 – Dependency-based formalism provides a way of expressing semantic re-  
182 lations that will constitute a good base for tackling the next steps in the  
183 analysis-chain, such as verb valence and thematic role studies (Agirre et  
184 al. 2006).
- 185 – The nature of the computational tools we have used for the preprocessing  
186 of the corpus to be tagged facilitates the establishment of dependency  
187 relations.
- 188 – The rich information involved when using the dependency model would  
189 allow transformation from trees to other means of representation.
- 190 – From our viewpoint, it is more straightforward to evaluate the relation  
191 between the elements that make up a sentence than the relation between  
192 elements included in parentheses, since the latter involves the additional  
193 task of determining where the parentheses start and end.
- 194 – In our opinion, dependency-based formalism is a more accurate method  
195 for annotating empty elements, such as *pro*, long-distance dependencies,  
196 and discontinuous constructions.

### 197 3.4. Theoretical and methodological basis

198 The design of a methodology for creating a treebank has different objectives:  
199 (i) to demonstrate the varieties of syntactic patterns of a language exhaus-  
200 tively; (ii) to remain correlated with the latest linguistic theories; (iii) to create  
201 an annotation scheme that can be used extensively in later research activities  
202 and computer-assisted practical solutions.

203 According to Hinrichs and Simov (2005), the relationship between the  
204 practice of treebank annotation and linguistic theorizing has become an im-  
205 portant subject of research. The advantage of assuming a particular theory is  
206 that it may solve many problems. The disadvantage, however, is that theories  
207 are unable to predict many aspects contained in the corpus. On the other  
208 hand, theory-neutral annotation schemes attempt to encode those grammat-  
209 ical properties that are distinguished by many, if not all, grammatical frame-  
210 works, without adhering to any particular linguistic theory. Theory-neutral  
211 annotations have the advantage of being more widely usable and of being less  
212 dependent on whatever version of a particular grammatical theory may have  
213 existed at the time when the treebank annotation scheme was determined.  
214 Since linguistic theories tend to change rapidly over time, and since treebank  
215 annotation is a labour-intensive and costly process, it is generally not feasi-  
216 ble to update Treebank annotations as a particular linguistic theory begins to  
217 change.

218 Proponents of the theory-dependent Treebank point out that the notion  
219 of a theory-neutral annotation is in itself an illusion, since any annotation  
220 scheme is the result of at least implicit linguistic theorizing. These scholars  
221 also point out that grounding an annotation scheme on a linguistic theory

222 tends to improve the consistency of the annotations, at least if the theory  
 223 provides explicit guidelines for the style of syntactic annotations. The Prague  
 224 Dependency Treebank (Hajic 1998) is a prime example of a theory-dependent  
 225 treebank.

226 With all these considerations in mind, and taking into account the literature  
 227 on tagging corpora in different languages, we decided to focus on certain pa-  
 228 rameters for determining the theoretical and methodological basis we needed  
 229 to build the Basque Treebank (BDT). The basic aspects addressed are as fol-  
 230 lows:

231 *3.4.1. Do we follow any theory?*

232 We follow Dependency Grammar model theory, so that for any dependency  
 233 relation we establish as head and dependent units the corresponding lexical  
 234 units. According to this model we annotate each word with its corresponding  
 235 dependency relation tag (cf. Figure 1). In annotating we have not taken into  
 236 account any specific linguistic theory; we determine the most appropriate  
 237 dependency relation tag, depending on the phenomena we have to deal with.  
 238 For instance, in the case of elided elements, such as *pro*, we adopt the gener-  
 239 ative approach in order to give a deeper linguistic analysis (the most important  
 240 features of our annotation practice are set out in section 4).

241 In addition, the analyses are motivated by a precise, comprehensive and  
 242 coherent theory of Basque grammar proposed by The Academy of the Basque  
 243 Language (Euskaltzaindia 1991 [1985], 1987, 1990, 1994, 1999).

244 *3.4.2. Which elements will be tagged?*

245 Our object of study is the sentence; i.e. the text enclosed between two full  
 246 stops (and also some other punctuation marks such as exclamation marks,  
 247 question marks and colons).

248 As well as the explicit elements making up the sentence, we have also  
 249 considered certain elided elements such as the *pro*.

250 In addition, long-distance dependencies and discontinuous constructions  
 251 are also annotated; that is, multiword lexical units (e.g. *bat egin* in (1)), name-  
 252 entities (e.g. *Henriette Aire* in (2)) and complex postpositions (e.g.. *kartelen*  
 253 *artetik* in (3)), obtained by IXATI as analysis units.

254 (1) *Proposamen-arekin bat egin zuen Espilondo-k*  
 Proposal-SC.COM To join AUX-PST-3SG-SG Espilondo-ERG  
 255 ‘Espilondo joins the proposal.’

256 (2) *Henrietta Aire-k olerki unibertsal-ari buruzko bere*  
 Henrietta Aire-ERG poetry universal-DAT about her-POS  
 257 *gogoeta-k azaldu-ko ditu*  
 thought-PL explain-FUT AUX-3PL-3SG  
 258 ‘Henriette Aire will explain her thoughts about universal poetry.’

- 259 (3) *Leiho-ko kartel-en arte-tik begirutzen du*  
 Window-GENLOC poster-GEN through- look AUX-3SG-3SG  
 260 ‘He/she looks through the posters affixed to the window.’

261 Furthermore, we have defined some auxiliary labels to tag units as multiword  
 262 when the previous IXATI module does not treat them as such (see 4.1).

### 263 3.4.3. *Which component will be the head in a dependency relation?*

264 The criteria for establishing dependency relations, and for distinguishing the  
 265 head and the dependent in such relations, are clearly of central importance  
 266 for dependency grammar. Such criteria have been discussed not only in the  
 267 dependency grammar tradition, but also within other frameworks where the  
 268 notion of syntactic head plays an important role, including all constituency-  
 269 based frameworks that subscribe to some version of *X* theory (Chomsky 1970;  
 270 Jackendoff 1977). Here are some of the criteria that have been proposed for  
 271 identifying a syntactic relation between a head (H) and a dependent (D) in a  
 272 construction (C) (Zwicky 1985; Hudson 1990):

- 273 – H determines the syntactic category of C and can often replace C;
- 274 – H determines the semantic category of C; D gives semantic specification;
- 275 – H is obligatory; D may be optional;
- 276 – H selects D and determines whether D is obligatory or optional;
- 277 – The form of D depends on H (agreement or government);
- 278 – The linear position of D is specified with reference to H.

279 It is clear that this list contains a mix of different criteria, some syntactic and  
 280 some semantic, and one may ask whether there is a single coherent notion  
 281 of dependency corresponding to all of them. Taking into account some of  
 282 these criteria, we take morphological and syntactic features, such as POS,  
 283 position of Dependent and Head, Agreement, and so on to link the head and  
 284 its corresponding dependent. The linguistic principles followed can be found  
 285 in (Aranzabe 2008).

286 For instance, in the case of noun phrases (NP) and prepositional phrases  
 287 (PP) the noun will be the head of such structure. In this approach we differ  
 288 from the generativists who consider the determiner as the head of the NPs  
 289 and the postposition as the head of the PPs.

290 In summary, it was decided to take lexical units as the basis of any depen-  
 291 dency relation for the following reasons:

- 292 – The previously developed syntactic approach was the dependency-oriented  
 293 surface syntax.
- 294 – In the next step we will address semantics. In our opinion, considering  
 295 lexical units as syntactic heads will fit better with our semantic work.



#### 296 3.4.4. The annotation scheme employed

297 In order to define the tagging system we have assumed the hierarchy proposed  
 298 in Carroll et al. (1998). They propose an annotation scheme in which each  
 299 sentence in the corpus is marked up with a set of grammatical relations,  
 300 specifying the syntactic dependency which holds between each head and its  
 301 dependent(s). Following this approach we have developed a tagset based on  
 302 hierarchies of grammatical relations (see Figure 1). In this paper we will  
 303 explain the use of this tagset.

### 304 4. The syntactic annotation

305 In this section we will present, using examples, the most significant features  
 306 of the manual annotation process. Before explaining the use of the tagset,  
 307 we will speak about the way we have annotated multiword expressions not  
 308 identified by the previous computational process. We will go on to explain  
 309 the annotation of a noun, heads and their dependents in a clause; and we will  
 310 then give an overview of annotation in subordinate clauses.

311 Next, we will show the annotation in constructions such as appositions,  
 312 predicative clauses, coordination and we will give some examples of the  
 313 empty category *pro*.

314 Finally, we will present Abar-Hitz (Díaz de Ilarraza et al. 2004), the appli-  
 315 cation implemented to help the annotators in their work.

#### 316 4.1. Discontinuous constituents: Multiword expressions

317 Multiword expressions are a problematic issue both in NLP tasks and in theo-  
 318 retic studies. First of all, there is no clear delimitation when defining different  
 319 types of multiword expressions; and even delimited, the list of multiwords  
 320 is not a complete and closed one, as is the case with simple words. Thus,  
 321 when tagging real corpora, it is quite common to find two or more simple  
 322 words that should be treated as a single unit; that is, words that the mod-  
 323 ule for the multiword treatment has not detected as such. This is the case  
 324 with multiword lexical units, entities, complex postpositions and complex  
 325 subordinating conjunctions.

326 With regard to multiword lexical units, entities and complex-postpositions,  
 327 we should say that, although our automatic processing offers good precision  
 328 and recall, some of them remain unidentified; and at the same time, no mul-  
 329 tiwords associated with subordinating conjunctions are detected, since there  
 330 is no specific treatment for them (e.g. *hori egin bitartean* /'while doing that':  
 331 a verb with a complex conjunction).

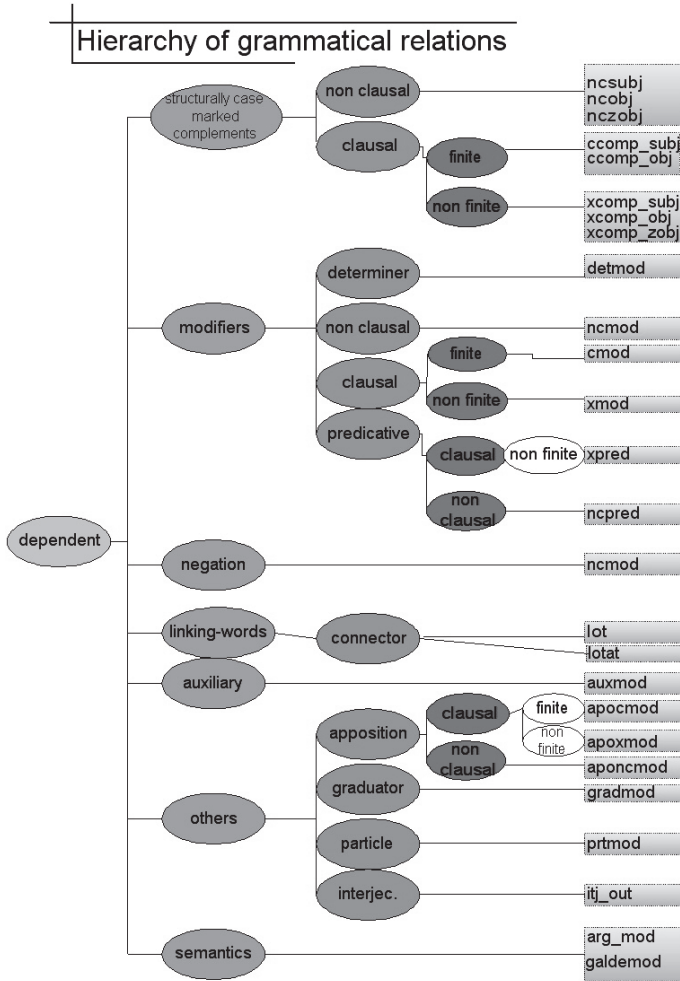


Figure 1. *Hierarchy of grammatical relations*

332 It is therefore necessary to define some auxiliary tags in order to analyze  
 333 these kinds of multiwords as a syntactic unit. We distinguish four auxiliary  
 334 tags for the four multiword expressions we are interested in:

- 335 **Haos** : components of multiword lexical units (e.g. *Ildo honetatik* /‘thus’:  
 336 a linking word).  
 337 **Menos** : components of multiword that express subordinating conjunctions.  
 338 (e.g. *hori egin bitartean* /‘while doing that’: a verb with a complex  
 339 conjunction).

340 **Entios**: component of a multiword entity. (e.g. *Peru Badiola Salazar*:  
 341 a proper name with its first and second surname).

342 **Postos**: component of a multiword postposition. (e.g. *liburuari dagokionez*/  
 343 ‘regarding the book’: a common name with a complex postposi-  
 344 tion).

345 In all of them, the annotator tags the component from left to the right. The  
 346 last component of the multiword should have the corresponding common  
 347 tag (e.g. *bitartean* is tagged as *xmod*, a non-finite subordinate clause). The  
 348 process is the same with discontinuous components (e.g. *baita ni ere* ‘me  
 349 too’, where *ere* is tagged as *lot*, a connector, and *baita* is tagged as *haos*).

#### 350 4.2. Noun heads and their dependents inside clauses

351 Dependency grammatical relations corresponding to non-clauses can be de-  
 352 scribed from two perspectives: i) as head of the relations (*ncsubj*, *ncobj*,  
 353 *nczobj*, *ncmod*, *ncpred* and *itj\_out*), and ii) as modifiers of heads (*detmod*,  
 354 *ncmod*, *aponcmo*d and *gradmod*).

355 When dealing with the structure of the non-clauses, we should say that we  
 356 are not concerned with understanding the internal structure of noun phrases.  
 357 We attempt to treat the phrases that are not clauses in a homogenous way.

358 Our approach is intended to provide consistent argument labelling that will  
 359 facilitate the automatic extraction of relational data, without attempting to  
 360 justify any theory.

##### 361 4.2.1. Head of the non-clause and the tagging representation

362 Basque is what is known as a head-final language, since heads tend to be  
 363 placed at the right-hand end of phrases. If we look at the structure of phrases  
 364 in Basque, we can see that the morphological marker is placed in the last  
 365 component of the phrase that carries it, regardless of the POS. Thus, the case  
 366 marker can be attached to the head<sup>4</sup> of a noun-clause as in (4) (e.g. *zalantza-k*)  
 367 or to a modifier of the head as in (5) (e.g. the adjective *altu-k*) and sometimes  
 368 to the determiner as in (6) (e.g. *hori-ek*):

369 (4) *Zenbait zalantza-k ezusteko bide-tik lortu*  
 Some doubt-ERG unexpected way-SG-ABL solve  
 370 *zuten argi-a*  
 AUX-PST-3SG-3PL light-SG-ABS  
 371 ‘Some doubts were solved in an unexpected way.’

372 (5) *Edozein mutil altu-k egiten du*  
 Any box tall-ERG do-IPFV AUX-PRS-3SG-3SG  
 373 ‘Any tall boy does it.’

374 (6) *Zalantza hori-ek ezusteko bide-tik lortu*  
 Doubt those-ERG unexpected way-SG-ABL solve  
 375 *zuten argi-a*  
 AUX-PST-3PL-3SG light-SG-ABS  
 376 ‘Those doubts were solved in an unexpected way.’

377 In order to maintain coherence in each relation when the element carrying the  
 378 declension-case/determiner and the noun head are not coincident, we decide  
 379 to include both elements<sup>5</sup> together explicitly in the description of the relation.  
 380 We consequently use a list of tuples to represent head/modifier relations in  
 381 the dependency tree. For example, a structurally case-marked complement  
 382 in which the complement is *nc* (non-clausal) has the following format:

- 383 – Case: the case marker by means of which the relation is established between  
 384 the head and the head of the phrase.
- 385 – Head: the governor of phrase.
- 386 – Head dependent.
- 387 – Case marker: the component of the phrase that carries the case.
- 388 – Syntactic function: the syntactic label assigned to the relationship.

389 The analyses of the phrases included in the following sentences exemplify  
 390 this formalization. In the phrase *zenbait zalantzak* in (4), *zalantzak* is the  
 391 element that carries the case marker and, at the same time, it constitutes the  
 392 head of phrase, so, the subject relation looks like the *ncsubj* dependency  
 393 shown below.

394 *detmod (- , zalantzak, zenbait)*  
 395 *ncsubj (erg, lortu, zalantzak, zalantzak, subj)*

396 In (6), the phrase *zalantza horiek*, *zalantza* is the head of the phrase, and so  
 397 we would add the component that carries the case marker, namely *horiek*.  
 398 Some of the relations associated to the NP follow:

399 *ncsubj (erg, lortu, zalantza, horiek, subj)*  
 400 *detmod (- , zalantza, horiek)*

#### 401 4.2.2. The dependency tags used

402 Regarding non-clause heads, we will distinguish two perspectives used to tag  
 403 phrases: i) the relations established between the noun and the verb: *ncsubj*,  
 404 *ncobj*, *nczobj*, *ncmod*, *ncpred* and *itj\_out* and ii) the modifiers of noun  
 405 heads: *detmod*, *ncmod*, *aponcmod* and *gradmod*.

406 Below we present the analysis and dependency-tree of the examples given  
 407 using the aforementioned dependency tags. The description of each of the  
 408 grammatical relations is extremely important, since it determines the number

409 and type of arguments needed for each relation (number of slots, the charac-  
 410 teristics of each one, etc.). This work will be very useful for future treatments,  
 411 for example in getting all this information into XML<sup>6</sup> format.

412 The description of all tags is presented in the appendix. We will give some  
 413 of them here to provide a better understanding of the tagging format used in  
 414 the dependencies.

415 Let us begin by giving the description of non-clausal tags. Some of them  
 416 have 5 slots (ncsubj, ncobj, nczobj, ncm<sup>7</sup>, ncpred), others 4 (ncmod<sup>8</sup>,  
 417 aponcm<sup>8</sup>, itj<sub>out</sub>), and others 3 (detmod, gradmod). Below we present  
 418 some examples of the representation of these relations:

- 419 ncsubj (Case, VerbHead, Head of NP, Case-marked element within NP, Role)
- 420 ncm<sup>7</sup> (Case, VerbHead, Head of NP, Case-marked element within NP)
- 421 ncpred (-, VerbHead, Head of NP, Case-marked element within NP)
- 422 ncm<sup>8</sup> (-, Noun Head, Case-marked element within NP)
- 423 detmod (-, Noun Head, Determiner)

424 In the examples bellow, the above-mentioned tags are used when tagging  
 425 the EPEC corpus. In this way it will better understand our annotation. For  
 426 each example, we will present the sentence in Basque and the translation  
 427 to English in two ways: English words in the same order as they are in the  
 428 Basque source text and the correct translation. Furthermore, the complete  
 429 relation set is added together with a graphical representation of the analysis  
 430 tree.

431 A characteristic in (7) is that the elements of the phrase linked to the verb  
 432 contain the case marker.

- 433 (7) *Zu-k galdu zenion beldurr-a itsaso-ari*  
 434 You-ERG lost AUX-PST-2SG-3SG-3SG fear-SG-ABS sea-DAT  
 435 *txiki-txiki-tatik*  
 436 childhood-ABL  
 437 ‘You have lost your fear of the sea since your childhood.’

438 Below we present the list of relations used for tagging the sentence:

- 437 ncsubj (erg, *galdu, zuk, zuk*, subj)
- 438 auxmod (-, *galdu, zenion*)
- 439 ncobj (abs, *galdu, beldurra, beldurra*, obj)
- 440 nczobj (dat, *galdu, itsasoari, itsasoari*, zobj)
- 441 ncm<sup>9</sup> (abl, *galdu, txiki-txikitatik, txiki-txikitatik, adlg*)

442 The dependency tree of this example can be seen in Figure 2.

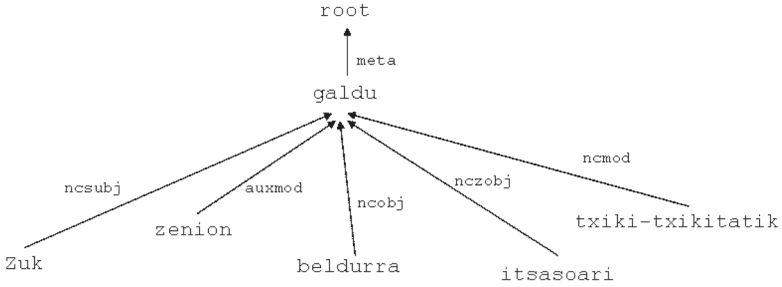


Figure 2. *Dependency tree for Zuk galdu zenion beldurra itsasoari txiki-txikitatik*

443 In (8) (see dependency tree in Figure 3), the noun *iritzi* is linked to the verb by  
 444 means of a *ncsubj* dependency relation although the case marker is included  
 445 in the determiner *hau* ‘this’ that modifies the noun. In this approach we make  
 446 no distinction between the predicative noun and verb; this is why in this  
 447 example, the noun *fruitu* is linked to the verb rather than to the noun *iritzi*.

- 448 (8) *Iritzi hau natura-ren behaketa zuzen-aren*  
 Opinion this-ABS nature-SG-GEN observation direct-GEN  
 449 *fruitu zen*  
 fruit-ABS be-PST-3SG  
 450 ‘This opinion was fruit of a direct observation of nature.’

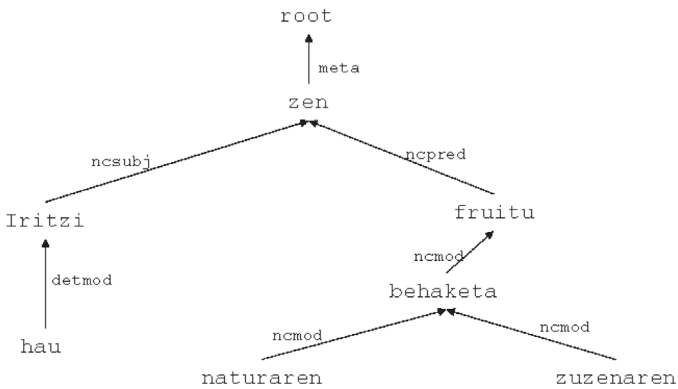


Figure 3. *Dependency tree for Iritzi hau naturaren behaketa zuzenaren fruitu zen*

451 In (9) (Figure 4) the *itj\_out* relation is illustrated. This relation differs from  
 452 the others insofar as it does not represent a common function in the sentence  
 453 structure, because it is a vocative or exclamation related to the direct style,  
 454 but it has been included in this group because it relates to a noun, *Valentine*,  
 455 and a to a verb, *bustitzen*.

456 (9) *Euri-ak ez zaitu bustitzen Valentine*  
 Rain-SG-ERG not AUX-PRS-2SG-3SG wet-IPFV Valentine-VOC  
 457 ‘The rain is not wetting you, Valentine.’

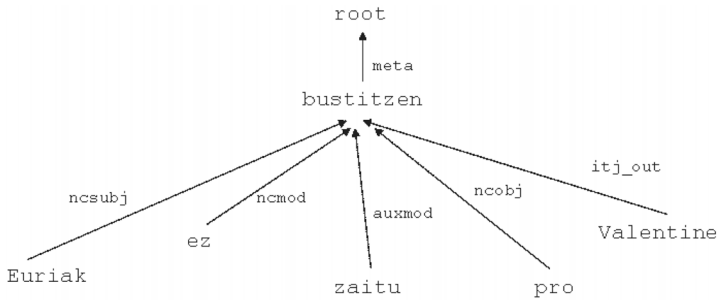


Figure 4. *Dependency tree for Euriak ez zaitu bustitzen, Valentine*

458 In (10) (see Figure 5) the internal relations of NP are shown; that is, the  
 459 dependents of the non-clausal head. Some types of NP structures have been  
 460 included in order to show their internal dependency relations. *Arrasateko* is  
 461 a noun modifier and the demonstrative *hau* appears to the right of the noun  
 462 while the quantifier *zenbait* and the ordinal *bigarren* precede the noun.  
 463 They are both linked to the noun.

464 (10) *Arrasate-ko zenbait familia-k bigarren tarifa hau*  
 Arrasate-GENLOC some famili-ERG second rate this-ABS  
 465 *kontratatu zuen*  
 hired AUX-3SG-3SG  
 466 ‘Some families from Arrasate hired this second rate.’

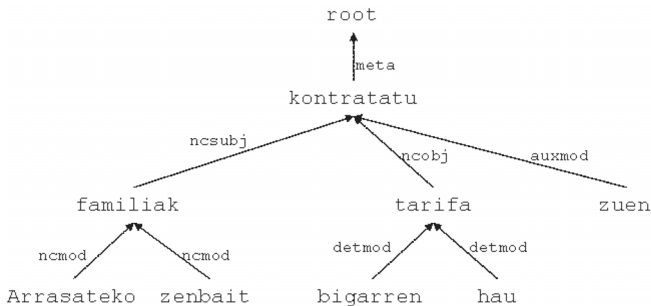


Figure 5. *The dependency tree for the sentence Arrasateko zenbait familiak bigarren tarifa hau kontratatu zuen*

467 In (11) (see Figure 6) we have the apposition structure classified in the hier-  
 468 archy in Figure 1 as others. It represents the relation between a noun and the  
 469 head of the preceding NP. In that case it is the relation between the heads of  
 470 two phrases. In the modifier relation expressed by *aponcmmod* the modifier is  
 471 *idazle* and the head *Axularrek*.

472 (11) *Axularrek, gure idazle handiak idatzi zuen*  
 Axular-ERG our writer great-SG-ERG write AUX-PST-3SG-3SG  
 473 *liburu hori*  
 book that-ABS  
 474 ‘Axular, our great writer, wrote that book.’

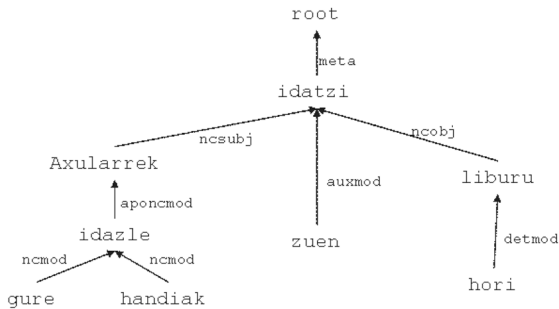


Figure 6. *The dependency tree for the sentence Axularrek, gure idazle handiak, idatzi zuen liburu hori*

475 Once we have revised, by means of examples, the tagset used for noun heads  
 476 and their dependents inside clauses, we will explain in a similar way the  
 477 relations defined for subordinate clauses.

#### 478 4.3. Subordinate clauses

479 Subordinate clauses are divided into complement and modifier. Sentence 12  
 480 (Figure 7) exemplifies the case in which the verb of the subordinate clause  
 481 (i.e. *dituela*) is finite. In this case, the verb of the subordinate clause is tagged  
 482 as *ccomp* and, depending on the function it performs with respect to the main  
 483 verb we will use *ccomp\_subj* or *ccomp\_obj* (see structure of the relations in  
 484 appendices). From here to the end of the paper, we will show the relations  
 485 for each example by means of their tree representation.

486 (12) *Gero, diote Euskal Herria-k zazpi probintzi*  
 Then say-PRS-3SG-3PL Euskal Herria-ERG seven province-ABS  
 487 *ditu-ela*  
 have-PRS-3SG-3PL-COMPL  
 488 ‘Then they say that the Basque Country has seven provinces.’



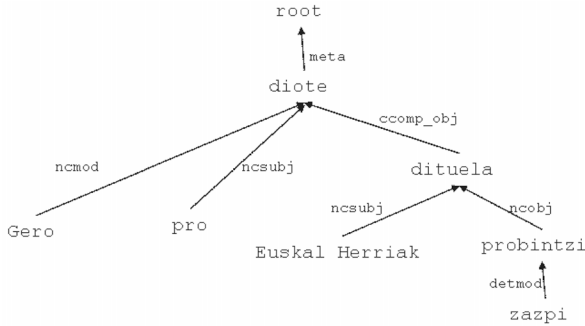


Figure 7. *The dependency tree for the sentence Gero diote Euskal Herriak zazpi probintzi dituela*

489 In Figure 8 we present a case of a sentence in which the verb of the subordinate  
 490 clause (i.e. *edateari* in (13)) is not finite (xcomp). Non-finite is represented  
 491 by the x in xcomp and, depending on the function it performs with respect  
 492 to the main verb we will use xcomp\_subj, xcomp\_obj, or xcomp\_zobj (see  
 493 structure of the relations in appendices).

- 494 (13) *Edateari eman nion*  
 495 drink-NMLZ-SG-DAT give AUX-PST-1SG-1SG  
 'I started drinking.'

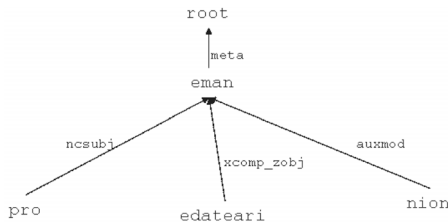


Figure 8. *The dependency tree for Edateari eman nion*

496 Below we show subordinate clauses that have the modifier function. Although  
 497 modifier subordinates can be of different types (time, cause, etc.), we use two  
 498 different variations of dependency tag, depending on the finiteness of the  
 499 verb of the subordinate clause; we have therefore associated the cmod tag  
 500 with finite verbs, (see (14)<sup>10</sup> in Figure 9) and xmod with non-finite verbs (see  
 501 (15)<sup>11</sup> in Figure 10). We have defined the following slots for the cmod relation:  
 502 i) clause type, ii) head of the main clause, iii) head of the subordinate clause,  
 503 and iv) auxiliary carrying the relational suffix. The xmod relation takes the  
 504 slots: i) clause type, ii) head of the main clause, iii) head of the subordinate  
 505 clause, and iv) word carrying the relational suffix.

506 In (14), we can also see a relative clause in which the *cmod* relation is  
 507 established between the verb *inguratzen* in the subordinate clause and the  
 508 noun *likidoari* of the main clause (antecedent).

509 (14) *Ezaugarri hau oso erraz froge daiteke*  
 property this-ABS very easily prove AUX-can-3SG  
 510 *koazerbatu-ak inguratzen ditu-en*  
 coacervate-PL-ABS surround-IPFV AUX-PRS-3SG-3SG-REL  
 511 *likido-ari koloratzaile desberdin-ak eranst*  
 liquid-SG-DAT colourant different-ABS-PL add-IPFV  
 512 *ba-dizkiogu*  
 COND-AUX-PRS-1PL-3SG-3PL  
 513 ‘This property can be proved very easily if we add different colour-  
 514 ants to the liquid surrounding the coacervate.’

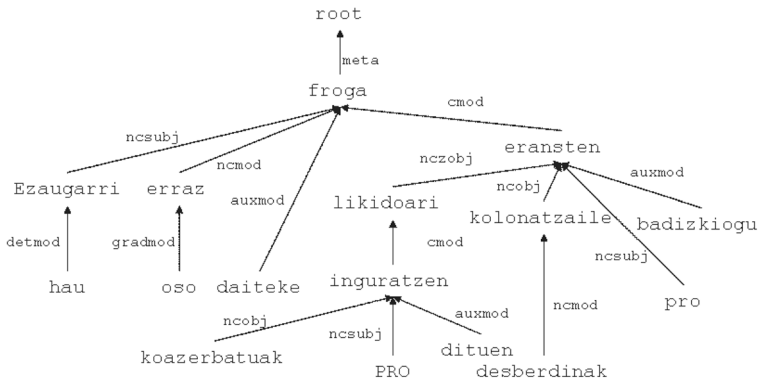


Figure 9. *The dependency tree for Ezaugarri hau oso erraz froge daiteke koazerbatuak inguratzen dituen likidoari koloratzaile desberdinak eranst badizkiogu*

515 In the rest of the cases, see example (15), the relation is given between two  
 516 verbs (main and subordinate clause).

517 (15) *Gertatu-tako-az jabetzen has-tean gaizki sentitu*  
 Happen-REL-INS realize-IPFV begin-TEMP terrible prove  
 518 *nintzen*  
 AUX-PST-1SG  
 519 ‘When I began to realise what had happened, I felt terrible.’

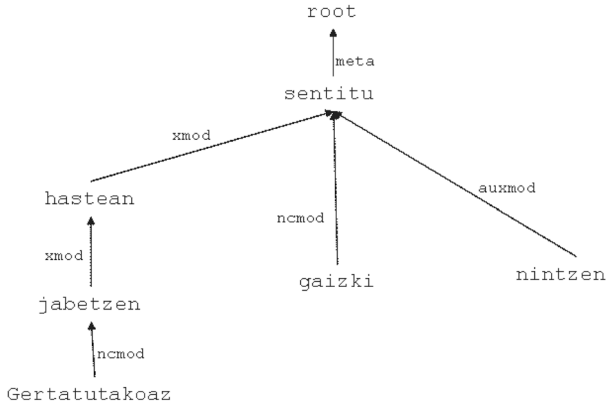


Figure 10. The dependency tree for Gertatutakoaz jabetzen hastean gaizki sentitu nintzen

520 4.4. Apposition and predicative clauses

521 In (16) (Figure 11) an apposition clause is shown. The apposition (re-  
 522 presented by the apo abbreviation, as in the non-clause tag) is an explanation  
 523 or specification of an element (either a complement or a modifier) that is the  
 524 head of the apposition. In the apposition clauses there are also two different  
 525 variations of apo dependency tag, depending on the finiteness of the verb of  
 526 the subordinate clause; we have therefore associated the apocmod tag with  
 527 finite verbs and apo<sub>x</sub>mod with non-finite verbs.

528 (16) *Jokin jokalari atzerriratu-a etorri da*  
 529 Jokin-ABS player abroad-ABS turn up AUX-PRS-3SG  
 ‘Jokin, the player who went abroad, has turned up.’

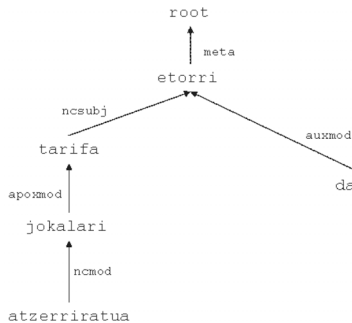


Figure 11. The dependency tree for Jokin, jokalari atzerriratu-a, etorri da

530 4.5. Analysis of coordination

531 Coordination is as problematic for Dependency Grammar formalism as for  
 532 other traditional theories. In order to capture the idea that the constituents  
 533 that are coordinated are at the same level, we have considered two options  
 534 extensively explained in the literature (Böhmová et al. 2003; Järvinen and  
 535 Tapanainen 1997): i) to presume one of the elements coordinated depends on  
 536 the other and ii) to add a new imaginary node, maintaining the coordinated  
 537 elements at the same level.

538 In our case, for computational reasons, we opt for the second one, which  
 539 is expressed by considering the coordinator element as a head of the coordi-  
 540 nate phrase; (17) (Figure 12) shows a case of noun phrase coordination that  
 541 illustrates our choice.

- 542 (17) *Horixe zen magoak eta nik*  
 543 *that-ABS be-PST-1SG illusionist-SG-ERG and I-ERG*  
*genuen sekretua*  
 544 *have-PST-1PL-3SG-REL secret-SG-ABS*  
 ‘That was the secret the illusionist and I had.’

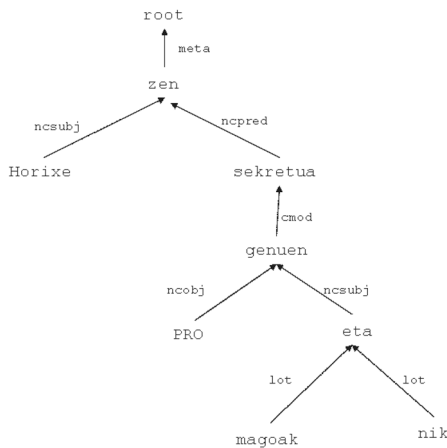


Figure 12. Example of the dependency tree of NP coordination

545 In the example, the coordinated elements *magoak* and *nik* are represented  
 546 at the same level and they have as their governor the connective *eta*, which  
 547 takes the dependency relation with respect to the verb, in this case *ncsubj*.

548 We use *emen* for copulative coordination, *aurk* for adversative, *haut* for  
 549 disjunctive, *espl* for explicative and so on.

550 The explanation given above could be extended to the coordination of more  
 551 than two elements.

552 4.6. The empty category *pro*

553 Basque displays a rich inflectional morphology. Indeed, it provides informa-  
 554 tion about the case (Absolute, Ergative or Dative) on either synthetic or  
 555 auxiliary predicates. Interestingly, it is possible for the argument phrase cor-  
 556 responding to one or several case markings not to appear in the sentence (the  
 557 so-called *pro*). However, precisely because the auxiliary displays case agree-  
 558 ment with this argument (which is a possibility with the so-called *pro-drop*  
 559 languages) we have assumed that this *pro* should be taken into account in the  
 560 sense that it belongs to the predicate when analyzing sentences. A subset of  
 561 50,000 words of EPEC has been manually annotated taking into account the  
 562 empty category as shown in (18) (Figure 13).

- 563 (18) *Begietatik igarri nionan ez*  
 564 Eye-SG-ABL figure out AUX-PST-1SG-3SG not  
 565 *zela bizi*  
 566 AZX.OST.3SG-COMPL live  
 567 ‘From his eyes I could see that he was dead.’

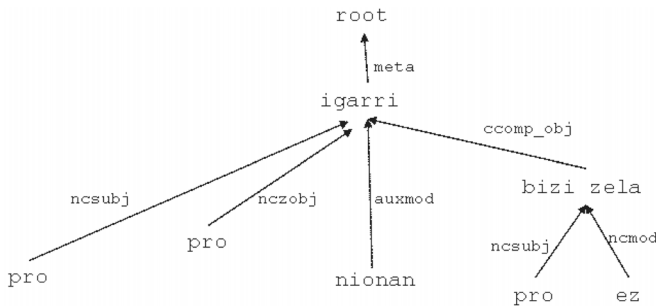


Figure 13. The dependency tree for *Begietatik igarri nionan ez zela bizi*

566 Once we have presented the details of the syntactic annotation we will briefly  
 567 explain the annotation tool used in the manual annotation process. We have  
 568 designed and implemented Abar-Hitz (Díaz de Ilarraza et al. 2004), a general  
 569 and friendly tool especially designed to help in the definition of dependen-  
 570 cies among the words of a sentence. It is important to emphasize that the  
 571 design of Abar-Hitz follows the general annotation schema we established  
 572 for representing linguistic information. It is part of a general environment  
 573 we have developed thus far, in which general processors and resources have  
 574 been integrated.

575 Abar-Hitz communicates with the user by means of a user-friendly interface  
 576 providing the following facilities:

- 577 1) It displays the morphosyntactic information obtained thus far, which has  
 578 previously been manually disambiguated in our specific tagging process.  
 579 The linguist is hardly requested to take this information into account when  
 580 beginning the syntactic tagging of a sentence. The tool is able to simul-  
 581 taneously use outputs from several tools (a morphological parser, a POS  
 582 tagger and a syntactic parser) to guide the annotator's decisions.
- 583 2) It graphically presents the dependency tree for each sentence. In addition,  
 584 the tree drawn can be graphically manipulated in such a way that the user  
 585 can change the tags and their fields, roll up subtrees, remove/add nodes,  
 586 remove/add connectors (dependencies) and so on. The changes in the tree  
 587 will be automatically verified when it is explicitly required or when the  
 588 window is closed.
- 589 3) It provides an environment for syntactic checking while tagging. We have  
 590 to take into account that mistakes can be made while tagging, both in the  
 591 number and type of slots and the name of the tag itself. Abar-Hitz avoids  
 592 these mistakes by showing specific pop-up menus where the only thing  
 593 the linguist can do is to select the appropriate tag.
- 594 4) It keeps track of unfinished sentences making it clear when these appear  
 595 on the screen.

596 Finally, the table below gives some figures for the occurrences and percent-  
 597 ages of the main dependency tags identified in EPEC:

Table 1. *Occurrences and percentages of dependency tags*

Dependency tag	Number	Percentage
ncmod	47817	34.17%
lot	18769	13.46%
auxmod	15172	10.61%
ncsubj	15287	10.73%
ncobj	11633	6.18%
detmod	7842	5.65%
xmod	5728	3.96%
xmod	4101	2.80%
ncpred	3548	2.50%
ccomp_obj	2029	1.42%
others	9361	8.53%

## 598 **5. Conclusions**

599 This paper has described the first formalization for the annotation of Basque  
 600 syntax using the Dependency Grammar Theory. We have started by setting  
 601 out the reasons for creating the BDT Treebank; i.e., a syntactically tagged

602 corpus. After considering and analyzing the principal possibilities that exist,  
 603 we decided to follow the formalism based on dependency relations, basically  
 604 for two reasons: first, because it is known to be more suitable for languages  
 605 with a free word order, like Basque; and second, because, apart from being  
 606 intuitive and easy to use, its flexibility allows new types of tags to be  
 607 introduced, such as those corresponding to thematic roles. This will be an  
 608 important aspect for any research we carry out in the future.

609 We have taken the step of analyzing the syntactic structures by explicitly  
 610 expressing the relation between the head and the dependent.

611 Additionally, we have found solutions to problems that have emerged when  
 612 describing some syntactic phenomena such as coordination, discontinuous  
 613 constituents, and so on. To date, 300,000 words have been annotated. The  
 614 Abar-Hitz annotation tool has been used in the annotation process. It was  
 615 created taking into account the characteristics of our XML linguistic anno-  
 616 tation.

617 To conclude, we would like to stress the urgent need for a syntactically  
 618 tagged corpus, which would serve to evaluate and improve the parser for  
 619 Basque that we are developing in the group. Furthermore, it will also be  
 620 a key ingredient for syntactic studies from a theoretical point of view. The  
 621 Treebank can be used to verify our linguistic intuitions.

## 622 **Appendix**

623 *A) All the dependency tags (29) with their general representation, and the*  
 624 *meaning of the abbreviations within the tags*

625 *aponcmod: (null, head, head of the apposition phrase, element with a declen-*  
 626 *sion case)*

627 *apocmod: (null, head, head of the apposition phrase, element with a subordi-*  
 628 *nating conjunction)*

629 *apoxmod: (null, head, head of the apposition, element with a subordinating*  
 630 *conjunction)*

631 *auxmod: (null, head, auxiliary)*

632 *ccomp\_subj: (comp/indirect style, head, head of the dependent, element with*  
 633 *a subordinating conjunction)*

634 *ccomp\_obj: (comp/indirect style, head, head of the dependent, element with*  
 635 *a subordinating conjunction)*

636 *cmmod: (relation, head, head of the dependent, element with a subordinating*  
 637 *conjunction)*

638 *detmod: (null, head of the phrase, determiner)*

639 *entios: (null, right-hand entity component, entity component)*

640 *galdemod: (null, head, reinforcing element)*

641 *gradmod: (null, head, graduator)*

642 *haos: (null, right-hand multiword component, multiword component)*

- 643 itj\_out: (null, head, head of the interjection, element with a declension case)  
 644 lot: (relation, conjunction, head)  
 645 lotat: (null, root, connector)  
 646 menos: (null, subordinating conjunction component, subordinating conjunction component)  
 647 ncmmod: (declension case\*, head \*, dependent\*, dependent\*)  
 648     • If it is a noun: (declension case, head, head of the phrase, element with a declension case)  
 649     • If it is the negation particle: (neg, head, ez, ez)  
 650     • If it is a complex postposition phrase: (the case of the complex postposition, head, postposition, postposition)  
 651     • If it is an adverb: (null, head, adverb, adverb)  
 652     • If it is an adjective modifying a noun: (null, head of the phrase, adjective, adjective)  
 653     • If it is the left-hand component of a compound: (null, head, component of the compound, component of the compound)  
 654 ncpred: (abs/pro, head, head of the phrase, element with a declension case)  
 655 ncsbj: (erg/abs/par, head, head of the phrase, element with a declension case, subj)  
 656 ncbj: (abs/par, head, head of the phrase, element with a declension case, obj)  
 657 nczobj: (dat, head, head of the phrase, element with a declension case, zobj)  
 658 postos: (kasua, right-hand postposition component, postposition component)  
 659 prtmod: (null, head, particle)  
 660 xcomp\_subj: (konp/zhg, head, element with a subordinating conjunction, element with a subordinating conjunction)  
 661 xcomp\_obj: (konpl/zhg, head, element with a subordinating conjunction, element with a subordinating conjunction)  
 662 xcomp\_zobj: (konpl, head, element with a subordinating conjunction, element with a subordinating conjunction)  
 663 xmod: (relation head, element with a subordinating conjunction, element with a subordinating conjunction)  
 664 xpred: (null, head, element with a subordinating conjunction, element with a subordinating conjunction)  
 665 arg\_mod: (semantic role)

678 *B) Meaning of the abbreviations within the tags*

<b>apo</b> apposition	<b>aux</b> auxiliary
<b>c</b> finite clause	<b>comp</b> complement
<b>enti</b> entity	<b>galde</b> reinforcing element
<b>grad</b> graduator	<b>ha</b> multiword
<b>itj</b> interjection	<b>lot</b> conjunction



<b>lotat</b> connector	<b>mod</b> modifier
<b>null</b> empty	<b>nc</b> non clause
<b>obj</b> object	<b>os</b> component
<b>out</b> element out of the clause	<b>post</b> postposition
<b>pred</b> predicative	<b>prt</b> particle
<b>subj</b> subject	<b>x</b> non finite clause
<b>zobj</b> indirect object	

## 679 **Bionotes**

680 Izaskun Aldezabal is Assistant Lecturer of Basque for Especial Purposes at the  
 681 University of the Basque Country. She received her PhD in Basque Philology  
 682 in 2004 at the University of the Basque Country. She is a researcher in the field  
 683 of Natural Language Processing. Her research interests include knowledge  
 684 bases, syntax and semantic of verbs. E-mail: izaskun.aldezabal@ehu.es

685 Maria Jesus Aranzabe is Assistant Lecturer of Basque for Especial Pur-  
 686 poses at the University of the Basque Country. She received her PhD in  
 687 Basque Philology in 2008 at the University of the Basque Country. She is  
 688 a researcher in the field of Natural Language Processing. Her research in-  
 689 terests include Treebank construction and shallow and deep syntax. E-mail:  
 690 maxux.aranzabe@ehu.es

691 Jose Mari Arriola is Assistant Lecturer of Basque for Especial Purposes at  
 692 the University of the Basque Country. He received his Master in Computa-  
 693 tional Linguistics from the University of Barcelona in 1991 and his PhD in  
 694 Basque Philology in 2000 at the University of the Basque Country. He is a  
 695 researcher in the field of Natural Language Processing. His research interests  
 696 include shallow and deep syntax. E-mail: josemaria.arriola@ehu.es

697 Arantza Diaz de Ilarraza is Aggregate Professor in Computer Languages  
 698 and Systems at the University of the Basque Country. She received her PhD in  
 699 Computer Sciences in 1990 at the University of the Basque Country. She is a  
 700 researcher in the field of Natural Language Processing. Her research interests  
 701 include development of Natural Language Processing Resources, Machine  
 702 Translation and Linguistic Annotations. E-mail: a.diazdeilarraza@ehu.es

## 703 **Notes**

704 \* This work was partly supported by the Spanish Ministry of Education and Science (IMLT:  
 705 General Model for Integration of Linguistic Tools and Resources: a proposal based on  
 706 XML standards. TIN2007-63173) and the Local Government of the Basque Country  
 707 (AnHITZ 2006: Language Technologies for Multilingual Interaction in Intelligent Envi-  
 708 ronments, IE06-185).

709 1. *pro*: elided syntactic arguments that typically arise when the predicate displays agreement  
 710 with the elided argument *pro* itself.

- 711 2. Example taken from Carroll et al. (1998).  
 712 3. [http://ufal.mff.cuni.cz/pcedt/doc/PCEDT\\_main.html](http://ufal.mff.cuni.cz/pcedt/doc/PCEDT_main.html).  
 713 4. Head is associated to any kind of analysis unit: multiwords or entities.  
 714 5. The decision, however, is not specific to Basque: more generally, it arises in the word-based  
 715 Constraint Grammar analyzer (Karlsson et al. 1995). Our manual tagging seeks to be as  
 716 compatible as possible with output obtained by the parser, for evaluation purposes. The  
 717 easiest way to achieve this involved adapting the original tagset as proposed by Carroll  
 718 et al. (1998), including, in some cases, an additional slot. Note that we do not change the  
 719 initial dependency philosophy; we merely adapt it to our needs.  
 720 6. XML stands for *Extended Generalized Markup Language*. This is the standard and gen-  
 721 eralized language used for tagging texts, namely, a metalanguage used for specifying sets  
 722 of tags as opposed to a single set of tags.  
 723 7. ncmmod represents the relation between the verb and the head of the non clausal phrase.  
 724 8. ncmobj represents the relation between the noun and the modifiers in the non clausal  
 725 phrase.  
 726 9. Verb modifier  
 727 10. Only the analysis of the subordinate clause is provided.  
 728 11. Same as in Example 14.

## 729 References

- 730 Abeillé, Anne (ed.)  
 731 2003 *Treebanks: Building and using parsed corpora*. Dordrecht, Boston & London:  
 732 Kluwer Academic Publishers.  
 733 Aduriz, Itziar  
 734 2000 *EUSMG: morfologiatiak syntaxira murriztapen gramatika erabiliz. Euskararen*  
 735 *desanbiguazio morfologikoaren tratamendua eta azterketa sintaktikoaren lehen*  
 736 *urratsak* [EUSMG: From morphology to syntax using Constraint Grammar.  
 737 Morphological disambiguation in Basque and first steps in syntax]. Donostia-  
 738 San Sebastián: University of Basque Country thesis.  
 739 Aduriz, Itziar, María Jesús Aranzabe, Jose Mari Arriola, Aitziber Atutxa, Arantza Díaz de  
 740 Illarraza, Nerea Ezeiza, Koldo Gojenola, Maite Oronoz, Aitor Soroa & Ruben  
 741 Urizar  
 742 2006 Methodology and steps towards the construction of EPEC, a corpus of written  
 743 Basque tagged at morphological and syntactic levels for automatic processing.  
 744 In Andrew Wilson, Paul Rayson & Dawn Archer (eds.), *Corpus linguistics*  
 745 *around the world*, 1–15. Netherlands: Rodopi.  
 746 Agirre, Eneko, Izaskun Aldezabal, Jone Etxeberria & Elixabete Pociello  
 747 2006 A preliminary study for building the Basque PropBank. Paper presented at  
 748 the International Conference on Language Resources and Evaluations, Genoa  
 749 (Italy), 22–28 May.  
 750 Aldezabal, Izaskun, Klara Ceberio, Itsaso Esparza, Ainhara Estarrona, Jone Etxeberria, Mikel  
 751 Iruskietia, Eli Izagirre & Larraitx Uria  
 752 2007a *EPEC (Euskararen Prozesamendurako Erreferentzia Corpusa) segmentazio-*  
 753 *mailan etiketatzeko eskuliburua* [A guide to tag EPEC (Reference Corpus for  
 754 the Processing of Basque) at segmentation level]. Technical Report, No. TR-11,  
 755 1–45. University of Basque Country (Spain).  
 756 Aldezabal, Izaskun, María Jesús Aranzabe, Jose Mari Arriola, Arantza Díaz de Ilarraza, Ain-  
 757 hara Estarrona, Enrique Fernandez, Uria Larraitx & Mikel Iruskietia  
 758 2007b *EPEC (Euskararen Prozesamendurako Erreferentzia Corpusa) dependentzie-*  
 759 *kin etiketatzeko eskuliburua* [A guide to tag EPEC (Reference Corpus for the

- 760 Processing of Basque) with dependences]. Technical Report, No. TR-12, 1–  
761 113. University of Basque Country (Spain).
- 762 Alegria, Iñaki, Artola Xabier & Kepa Sarasola  
763 1997 Improving a robust morphological analyser using lexical transducers. In Ruslan  
764 Mitkov and Nicolas Nicolov (eds.), *Recent advances in Natural Language*  
765 *Processing*, 97–110. Amsterdam & Philadelphia: John Benjamins.
- 766 Alegria, Iñaki, Olatz Arregi, Nerea Ezeiza & Izaskun Fernandez  
767 2006 Lessons from the development of a named entity recognizer. *Procesamiento*  
768 *del Lenguaje Natural* 36. 25–37.
- 769 Aranzabe, María Jesús  
770 2008 *Dependentzia-ereduan oinarritutako baliabide sintaktikoak: zuhaitz-bankua*  
771 *eta gramatika konputazionala* [Linguistic resources based on dependency for-  
772 malism: A treebank and a computational grammar]. Donostia-San Sebastián:  
773 University of Basque Country thesis.
- 774 Artiagoitia, Xabier  
775 2002 The functional structure of the Basque noun phrase. *Journal of Basque Lin-*  
776 *guistics and Philology* 44. 73–90.
- 777 Artola, Xabier, Arantza Díaz de Ilarraza, Nerea Ezeiza, Koldo Gojenola, Gorka Labaka, Aitor  
778 Sologaitoa & Aitor Soroa  
779 2005 A framework for representing and managing linguistic annotations based on  
780 typed feature structures. Paper presented at the Recent Advances on NLP  
781 (RANLP), Borovets (Bulgaria), 18–24 September.
- 782 Bengoetxea, Kepa & Koldo Gojenola  
783 2006 *Desarrollo de un analizador sintáctico estadístico basado en dependencias*  
784 *para el euskera* [Development of a statistical parser for Basque]. *Procesamiento*  
785 *del lenguaje natural* 39. 5–12.
- 786 Böhmová, Alena, Jan Hajic, Eva Hajicová & Barbora Hladká  
787 2003 The PDT: A 3-level annotation scenario. In Anne Abeillé (ed.), *Treebanks:*  
788 *Building and using parsed corpora*, 103–127. Dordrecht, Boston & London:  
789 Kluwer Academic Publishers.
- 790 Brants, Thorsten, Skut Wojciech & Hans Uszkoreit  
791 2003 Syntactic annotation of a German newspaper corpus. In Anne Abeillé (ed.),  
792 *Treebanks: building and using parsed corpora*, 73–87. Dordrecht, Boston &  
793 London: Kluwer Academic Publishers.
- 794 Bunt, Harry, John Carroll & Giorgio Satta Giorgio  
795 2004 New Developments in Parsing Technology. *Text, Speech, and Language Tech-*  
796 *nology* 32(3). 439–442.
- 797 Carroll, John, Ted Briscoe & Antonio Sanfilippo  
798 1998 Parser evaluation: A survey and a new proposal. Paper presented at the Inter-  
799 national Conference on Language Resources and Evaluations, University of  
800 Granada (Spain), May.
- 801 Charniak, Eugene  
802 2000 A maximum-entropy-inspired parser. Paper presented at the first conference of  
803 the North American chapter of the Association for Computational Linguistics,  
804 Seattle (Washington), 29 April–4 May.
- 805 Chomsky, Noam  
806 1970 Remarks on nominalizations. In Roman Jacobs & Peter Rosenbaum (eds.),  
807 *Readings in English Transformational Grammar*, 184–221. Massachusetts:  
808 Ginn Waltham.

- 809 Collins, Michael  
810 1997 Three generative lexicalised models for statistical parsing. Paper presented at  
811 the Eighth Conference of the European Chapter of the Association for Com-  
812 putational Linguistics, Madrid (Spain), 7–12 July.
- 813 Collins, Michael  
814 2000 Discriminative reranking for natural language parsing. Paper presented at the  
815 Seventeenth International Conference on Machine Learning, Stanford (Cali-  
816 fornia), 29 June–2 July.
- 817 Díaz de Ilarraza, Arantza, Garmendia Aitzpea & Maite Oronoz  
818 2004 Abar-Hitz: An annotation tool for the Basque Dependency Treebank. Paper pre-  
819 sented at the International Conference on Language Resources and Evaluation,  
820 Lisbon (Portugal), 24 May.
- 821 Eguzkitza, Andolin  
822 1993 Adnominals in the grammar of Basque. In José Ignacio Hualde & Jon Ortiz de  
823 Urbina (eds.), *Generative studies in Basque linguistics*, 163–187. Amsterdam  
824 & Philadelphia: John Benjamins.
- 825 Euskaltzaindia  
826 1991 [1985] *Euskal Gramatika: Lehen Urratsak – I* [A Grammar of Basque: The first steps –  
827 I]. Bilbo: Euskaltzaindia.
- 828 Euskaltzaindia  
829 1987 *Euskal Gramatika: Lehen Urratsak – II* [A Grammar of Basque: The first  
830 steps – II]. Bilbo: Euskaltzaindia.
- 831 Euskaltzaindia  
832 1990 *Euskal Gramatika: Lehen Urratsak – III (Lokailuak)* [A Grammar of Basque:  
833 The first steps – III (the connectives)]. Bilbo: Euskaltzaindia.
- 834 Euskaltzaindia  
835 1994 *Euskal Gramatika: Lehen Urratsak – IV (Juntagailuak)* [A Grammar of Basque:  
836 The first steps – IV (the conjunctions)]. Bilbo: Euskaltzaindia.
- 837 Euskaltzaindia  
838 1999 *Euskal Gramatika: Lehen Urratsak – V (Mendeko perpausak – I)* [A Grammar  
839 of Basque: the first steps – V (subordinated clauses – 1)]. Bilbo: Euskaltzaindia.
- 840 Ezeiza, Nerea, Itziar Aduriz, Iñaki Alegria, Jose Mari Arriola & Ruben Urizar  
841 1998 Combining stochastic and rule-based methods for disambiguation in agglutina-  
842 tive languages. Paper presented at the 36th Annual Meeting of the Association  
843 for Computational Linguistics and 17th International Conference on Compu-  
844 tational Linguistics, Montreal (Canada), 10–14 August.
- 845 Ezeiza, Nerea  
846 2003 *Corpusak ustiatzeko tresna linguistikoak. Euskararen etiketatzaile morfosin-  
847 taktiko sendo eta malgua* [Linguistic tools for exploiting corpora. The Basque  
848 Morphosyntactic robust and flexible tagger]. Donostia-San Sebastián: Univer-  
849 sity of Basque Country thesis.
- 850 Goenaga, Patxi  
851 1991 *Gramatika bideetan* [In the paths of grammar]. Donostia-San Sebastián: Erein.
- 852 Hajiè, Jan  
853 1998 Building a syntactically annotated corpus: The Prague Dependency Treebank.  
854 In Eva Hajičová (ed.), *Studies in honour of Jarmila Panevová*, 106–132. Prague:  
855 Karolinum Charles University Press.
- 856 Hinrichs, Erhard & Kiril Simov  
857 2005 Special issue on treebanks and linguistic theories. *Research on Language and  
858 Computation* 3(1).
- 859 Hudson, Richard  
860 1990 *Word Grammar*. Oxford: Basil Blackwell Publishers.

- 861 Jackendoff, Ray  
 862     1997     *The architecture of the language faculty*. Cambridge, MA: The MIT Press.
- 863 Järvinen, Timo & Pasi Tapanainen  
 864     1997     *A dependency parser for English*. Technical Report, No. TR-1, 1–43. University  
 865     of Helsinki (Finland).
- 866 Kakkonen, Tuomo  
 867     2005     Dependency treebanks: Methods, annotation schemes and tools. Paper pre-  
 868     sented at the 15<sup>th</sup> Nordic Conference of Computational Linguistics, Joensuu  
 869     (Finland), 21–22 May.
- 870 Karlsson, Fred, Atro Voutilainen, Juha Heikkilä & Arto Anttila  
 871     1995     *Constraint grammar: A language-independent system for parsing unrestricted*  
 872     *text*. Berlin & New York: Mouton de Gruyter.
- 873 Laka, Itziar  
 874     1993     Unergatives that assign ergative, unaccusatives that assign accusative. In Jona-  
 875     than Bobaljik & Colin Phillips (eds.), *Case and agreement* (Linguistics 18),  
 876     149–172. Cambridge, MA: The MIT Press.
- 877 Lin, Dekang  
 878     1995     A dependency-based method for evaluating broad-coverage parsers. Paper pre-  
 879     sented at the International Joint Conference on Artificial Intelligence (IJCAI-  
 880     95), Montreal (Canada), 19–20 August.
- 881 Marcus, Mitchell, Beatrice Santorini & Mary Ann Marcinkiewicz  
 882     1993     Building a large annotated corpus of English: The Penn Treebank. *Computa-  
 883     tional Linguistics* 19(2). 313–330.
- 884 Mel'cuk, Igor Aleksandrovic  
 885     1988     *Dependency syntax: Theory and practice*. Albany: State University of New  
 886     York Press.
- 887 Montemagni, Simonetta, Francesco Barsotti, Marco Battista, Nicoletta Calzolari, Ornella  
 888     Corazzari, Alessandro Lenci, Antonio Zampolli, Francesca Fanciulli, Maria  
 889     Massetani, Remo Raffaelli, Roberto Basili, Maria Teresa Pazienza, Dario Sara-  
 890     cino, Fabio Zanzotto, Mana Nana, Fabio Pianesi & Rodolfo Delmonte  
 891     2003     Building the Italian Syntactic-Semantic Treebank. In Anne Abeillé (ed.), *Tree-  
 892     banks: Building and using parsed corpora*, 189–210. Dordrecht, Boston &  
 893     London: Kluwer Academic Publishers.
- 894 Oflazer, Kemal, Zeynep Dilek & Gökhan Tür  
 895     1999     Design for a Turkish treebank. Paper presented at the Workshop on Linguisti-  
 896     cally Interpreted Corpora, Bergen (Norway).
- 897 Sampson, Geoffrey  
 898     2003     Thoughts on two decades of drawing trees. In Anne Abeillé (ed.), *Treebanks:  
 899     Building and using parsed corpora*, 23–41. Dordrecht, Boston & London:  
 900     Kluwer Academic Publishers.
- 901 Skut, Wojciech, Brigitte Krenn, Thorsten Brants & Hans Uszkoreit  
 902     1997     An annotation scheme for free word order languages. Paper presented at the  
 903     Fifth Conference on Applied Natural Language Processing (ANLP'97), Wash-  
 904     ington (DC), 31 March–3 April.
- 905 Sleator, Daniel & Davy Temperley  
 906     1993     Parsing English with a link grammar. Paper presented at the Third International  
 907     Workshop on Parsing Technologies, Tilburg, 13 August.
- 908 Tapanainen, Pasi & Atro Voutilainen  
 909     1994     Tagging accurately – don't guess if you know. In *Proceedings of Applied Natural  
 910     Language Processing (ANLP'94)*, University of Stuttgart, 13–15 October.

- 911 Tapanainen, Pasi & Timo Järvinen  
912 1997 A non-projective dependency parser. In *Proceedings of the 5th Conference on*  
913 *Applied Natural Language Processing*, Washington (DC), 31 March–3 April.  
914 Tapanainen, Pasi & Timo Järvinen  
915 1998 Dependency concordances. *International Journal of Lexicography* 11(3). 187–  
916 203.
- 917 Tesnière, Lucien  
918 1959 *Éléments de Syntaxe Structurale*. Paris: Librairie Klincksieck.
- 919 Trask, Robert Lawrence  
920 2003 The noun phrase: Nouns, determiners and modifiers; pronouns and names. In  
921 José Ignacio Hualde & Jon Ortiz de Urbina (eds.), *A Grammar of Basque*,  
922 113–170. Berlin & New York: Mouton de Gruyter.
- 923 Zabala, Igone  
924 2003 Nominal Predication. In José Ignacio Hualde & Jon Ortiz de Urbina (eds.), *A*  
925 *Grammar of Basque*, 426–446. Berlin & New York: Mouton de Gruyter.
- 926 Zwicky, Arnold M.  
927 1985 Heads. *Journal of Linguistics* 21. 1–29.