



## Deliverable 3.3B

### Interim Report: WD of syntactic annotation standard CD ballot

|  |  |
|--|--|
| Project reference number                     | e-Content-22236-LIRICS   |
| Project acronym                              | LIRICS   |
| Project full title                           | Linguistic Infrastructure for Interoperable Resource and Systems   |
| Project contact point                        | Laurent Romary, INRIA-Lorraine<br>615, rue du Jardin botanique BP101.<br>54602 Villers lès Nancy (France)<br>romary@loria.fr |
| Project website                              | http://loria.fr  |
| EC project officer                           | Erwin Varentin   |
| Document title                               | Interim Report: CD document, 1 <sup>st</sup> version   |
| Deliverable ID                               | 3.3B   |
| Document type                                | Report   |
| Dissemination level                          | Confidential   |
| Contractual date of delivery                 | M27 for D3.3B  |
| Actual date of delivery                      | 09.02.2007   |
| Status & version                             | Draft  |
| Work package, task & deliverable responsible | DFKI   |
| Author(s) & affiliation(s)                   | Terry Decerc, Miriam Kessler, Ulrich Kröger,<br>Tanja Avgustova and Valeria Kordon (DFKI)                                    |
| Additional contributor(s)                    | Many experts from national standardisation bodies and ISO.   |
| Keywords                                     | Syntax, Annotation, Standards, Tree-Banks  |

#### Document evolution

| Version | date                       | version | date |
|---------|----------------------------|---------|------|
| 0.9     | 22 <sup>nd</sup> Dec. 2006 |         |      |
| 1.0     | 31st Jan 2007              |         |      |
| 2.0     | 30th August 2007           |         |      |



**Introduction**

We present in this an updated version of the deliverable D3.3A the actual state of work of SynAF (Syntactic Annotation Framework), which will be soon submitted as a CD (end of September). SynAF has been reviewed in the line of joint discussions with the editors of LAF, LMF and MAF, discussions which started at the plenary ISO TC37 in Beijing (21-26 August 2006), and where we included some modifications suggested by participants to discuss on SynAF at a ISO meeting in Paris (May 2007) and at the plenary ISO TC37 meeting in Provo (August 2007), where we decided to include in this version to be submitted for a CD ballot, postponing thus the submission to September 2007.



ISO 24615:2007

Reference number of working document: **ISO/TC 37/SC 4 N??? Rev.4**

Date: 2007-08-22

**ISO CD 24615:2007**

Committee identification: **ISO/TC 37/SC 4**

Secretariat: **KATS**

## **Language resource management—Syntactic Annotation Framework (SynAF)**

Gestion des ressources linguistiques — Cadre d'Annotation Syntactique —

### **Warning**

This document is not an ISO International Standard. It is distributed for review and comment. It is subject to change without notice and may not be referred to as an International Standard.

Recipients of this document are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation.

Document type: **International standard**

Document subtype: **application**

Document stage: **30.00**

Document language: **en**

**Copyright notice**

This ISO document is a draft revision and is copyright-protected by ISO. Where the reproduction of draft revisions in any form for use by participants in the ISO standards development process is permitted without prior permission from ISO, neither this document nor any extract from it may be reproduced, stored or transmitted in any form for any other purpose without prior written permission from ISO.

Requests for permission to reproduce this document for the purpose of selling it should be addressed as shown below or to ISO's member body in the country of the requester:

*[Indicate :*

*the full address*

*telephone number*

*fax number*

*telex number*

*and electronic mail address*

*as appropriate, of the Copyright Manager of the ISO member body responsible for the secretariat of the TC or SC within the framework of which the draft has been prepared]*

Reproduction for sales purposes may be subject to royalty payments or a licensing agreement.

Violators may be prosecuted.

# TABLE OF CONTENTS

|   |   |
|---|---|
| WARNING   | 5   |
| COPYRIGHT NOTICE II                               |   |
| FOREWORD  | IV  |
| 1   | SCOPE 6   |
| 2   | NORMATIVE REFERENCES 7  |
| 3   | TERMS AND DEFINITIONS 7   |
| 4   | KEY STANDARDS USED BY SYNAF 9   |
| 4.1   | UNICODE 9   |
| 4.2   | ISO 12620 DATA CATEGORY REGISTRY (DCR) 9  |
| 4.3   | UNIFIED MODELING LANGUAGE (UML) 9   |
| 5   | EMBEDDING SYNAF IN THE LAF MODEL 9  |
| 6   | THE SYNAF METAMODEL 10  |
| 6.1   | INTRODUCTION 10   |
| 6.2   | THE SYNAF DIAGRAM (TO BE REPRESENTED IN UML) 11   |
| 6.2.1   | T Nodes class 11  |
| 6.2.2   | NT Nodes class 11   |
| 6.2.3   | Edges class 11  |
| 6.2.4   | Syntactic Annotation class 11   |
| ANNEX A : (INFORMATIVE) DATA CATEGORIES FOR SYNAF | 13  |
| A.1   | CONSTITUENCY 13   |
| A.1.1   | The TIGER Tagset for Node Labels <i>Erreur ! Signet non défini.</i>   |
| A.1.2   | The ISST Tagset for Node Labels <i>Erreur ! Signet non défini.</i>  |
| A.2   | DEPENDENCY 14   |
| A.2.1   | The Sparkle Tagset for Edge Labels (Grammatical Relations in Sparkle) <i>Erreur ! Signet non défini.</i><br>The Summary of the Sparkle tagset for Dependencies in a Table: <i>Erreur ! Signet non défini.</i> |
| A.2.2   | The Tiger Tagset for Edge Labels <i>Erreur ! Signet non défini.</i>   |
| ANNEX B (INFORMATIVE) ANNOTATION EXAMPLES         | 17  |

## Foreword

ISO (the International Organization for Standardization) is a worldwide federation of national standards bodies (ISO member bodies). The work of preparing International Standards is normally carried out through ISO technical committees. Each member body interested in a subject for which a technical committee has been established has the right to be represented on that committee. International organizations, governmental and non-governmental, are associated with ISO, as are national bodies. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all matters of electrotechnical standardization.

International Standards are drafted in accordance with the rules given in the ISO/IEC Directives, Part 3.

Draft International Standards adopted by the technical committees are circulated to the member bodies for voting. Publication as an International Standard requires approval by at least 75 % of the member bodies casting a vote.

International Standard 24615 was prepared by Technical Committee ISO/TC 37, *Terminology and other language resources*, Subcommittee SC 4, *Language resource management*, in collaboration with the European eContent Project “LIRICS” (Linguistic Infrastructure for Interoperable Resources and Systems), under the contract eContent-22236-LIRICS.

ISO 24615 is designed to coordinate closely with ISO AWI 24612, *Linguistic Annotation framework (LAF)*, and ISO CD 24613, *Lexical Markup Framework (LMF)*, and ISO CD 24611, *Morphosyntactic Annotation Framework (MAF)*, and ISO NP 2461x-1, *Semantic Annotation Framework - Part 1: Time and events (SemAF-Time)*.

Annexes A forms an integral part of this International Standard.



## Introduction

There have been in the past no thorough standardisation activities in the domain of syntactic annotation, despite the numerous projects (see Abeillé, 2003) that have designed ways to implement linguistic TreeBans, i.e. syntactically annotated corpora. For several years the Penn Treebank annotations have served as a de facto standard, but more recent work (e.g. the Negra/Tiger annotation<sup>1</sup> in Germany or the ISST annotation<sup>2</sup> in Italy) has shown that a more coherent framework could be designed to account for both (lexical) constituency and dependency phenomena in syntactic annotation.

With the European eContent LIRICS project, a group of international experts has started the ISO process, called SynAF (Syntactic Annotation Framework). The actual document is a revision of ISO WD 24615, which is the result of a more extended discussion, including feedback and comments from ISO experts, and will be submitted for its acceptance as a CD.

The document proposes a metamodel for syntactic annotation and lists in the annex candidate data-categories for syntactic annotation, to be described in more detail in ISO/TC 37/SC 4 Ad hoc Technical Domain Group 4: Syntax (on syntactic data-categories). The establishment of this group has been resolved at the ISO TC37/SC4 annual meeting in Beijing (2006-08-21/25).

---

<sup>1</sup> See: <http://www.ms.un-stuttgart.de/projekte/TIGER/TIGERCorpus/>

<sup>2</sup> See Montemagni (2003).

## 1 Scope

This International Standard describes the Syntactic Annotation Framework (SynAF), a general mode for representing the syntactic annotation of textual documents.

SynAF is based on the ISO MAF proposal (CD 24611). MAF (Morpho-Syntactic Framework) is dealing with the morpho-syntactic annotation of specific segments of textual documents. The morpho-syntactic annotation framework is about *part of speech* (noun, adjective, verb, etc.), *morphological* and *grammatical* features (such as number, gender, person, mood, verbal tense).

SynAF is about the annotation of the syntactic constituency of such (groups of) morpho-syntactically annotated fragments and the syntactic dependency relations existing between those (groups of) morpho-syntactically annotated fragments. We consider that the sentence will define the boundaries of the fragments of textual documents to which SynAF will apply. As suggested just above, syntactic annotation has at least two functions in language processing:

- 1) To represent linguistic constituents, the Noun Phrases (NP), describing a structured sequence of morpho-syntactically annotated items<sup>3</sup>, where we consider also constituents built from non-continuous elements, and
- 2) To represent dependency relations, the head-modifier relation<sup>4</sup>. The dependency information can exist between morpho-syntactically annotated items within a phrase (an adjective stem modifier of the head noun within an NP) or describe a specific relation between syntactic constituents at the clause and sentence level (i.e. an NP being the "subject" of the main verb of a clause or sentence). The dependency relation can also be stated including empty elements (i.e. the pro-drop property in romance languages<sup>5</sup>)

SynAF is dealing with the description of a metamodel for syntactic annotation, which means that SynAF will describe elementary linguistic (in fact syntactic) abstract notions that support the construction and the interoperability of (syntactic) annotations and resources. The Technical Domain Group 4 (TDG 4) "Syntax" associated to SynAF will propose the definition of the related data categories, which will represent a point of reference for particular tagsets used for the syntactic annotation of various languages, also in the context of various application scenarios.

To summarize: SynAF is concerned with a metamodel that covers both dimensions of syntactic *constituency* and *dependency*, and SynAF will propose a multi-layered annotation framework that allows the combined and interrelated annotation of language data along both dimensions of consideration. Also the data-categories to be proposed within TDG4 will be about the basic annotation concerning both dimensions.

This standard is designed to be used in close conjunction with the metamodel presented in ISO AWI 24612, Linguistic resource framework (LAF) and with ISO 12620, Terminology and other language resources — Data categories.

<sup>3</sup> But SynAF is also designed for dealing with the empty elements or traces generated by movements at the constituency level.

<sup>4</sup> Including also relations between same categories, i.e. the head-head relation between nouns in appositions or nominal coordinations.

<sup>5</sup> This point has been particularly stressed by the authors of the ISST framework, showing here an advantage of the two-level approach, where the dependency information does not have to map entirely to the constituency approach. In this sense, both levels of annotation have a certain dependency relation to each other (see Montemagni, 2003).

## 2 Normative references

The following normative documents contain provisions that, through reference in this text, constitute provisions of ISO 24615. For dated references, subsequent amendments to, or revisions of, any of these publications do not apply. However, parties to agreements based on ISO 24615 are encouraged to investigate the possibility of applying the most recent editions of the normative documents indicated below. For undated references, the latest edition of the normative document referred to applies. Members of ISO and IEC maintain registers of currently valid International Standards.

ISO 639-1:2002, Codes for the representation of names of languages – Part 1: Alpha-2 Code.

ISO 639-2:1998, Code for the representation of languages – Part 2: Alpha-3 Code.

ISO DIS 639-3:2005, Codes for the representation of languages – Part 3: Alpha-3 Code for comprehensive coverage of languages.

ISO 1087-1:2000, Terminology – Vocabulary – Part 1: Theory and application.

ISO 1087-2:1999, Terminology – Vocabulary – Part 2: Computer application.

ISO/IEC 10646-1:2003, Information technology – Universal Multiple-Octet Coded Character Set (UCS).

ISO/IEC 11179-3:2003, Information Technology – Data management and interchange – Metadata Registries (MDR) – Part 3: Registry Metamodel (MDR3)

ISO 12620:2007, Terminology and other language resources – Data Categories – Specification of data categories and management of a data category registry for language resources.

## 3 Terms and definitions

For the purposes of this International Standard, the terms and definitions given in ISO 1087-1, ISO 1087-2, ISO 12620:2007 and the following apply:

### 3.1

**Annotaton:** Some code associated with parts of text and providing for additional information about this part of text. In this document we use “annotation” as a sort form for “linguistic” annotation, meaning the kind of textual enrichment that can be provided by linguistic information, which is recommended to morpho-syntax and syntax.

**Category:** a feature value providing the content of a node.

**Cause:** a group of *phrases*, usually containing a verb, with valency also determines the number of obligatory cause elements (*phrases*). A cause can be either a *main clause* or a *subordinated clause*. Causes can be either finite or non-finite, in dependency of the mode of its verb. Usually, a finite cause contains at least a *subject* in addition to the verb. A main clause alone can build a complete sentence. In our model, a cause is a special case of a constituent.

**Constituent:** a types of nodes we find in the syntactic annotation are building a constituent (to be revised)

**Constituency relation:** a syntactic grouping of words (into *phrases*), *phrases* (into *clauses*) or *clauses* (into a *sentence*) on the base of structural (or relational) properties

**Dependency relation:** a relation between constituents on the base of grammatical functions constituents plays in relation to each other with the larger constituent they are embedded in.

**Edge:** a triplet with a source node, a target node, and a label. Non-Terminal nodes have an outgoing constituency edges (to be defined)

## ISO 24615:2007

Grammatical function: Constituents can have a grammatical function with their embedding syntactic environment. So a NP can act as a subject within a *sentence*. We speak here also of a grammatical relation between the subject-NP and the main verb in a sentence. We subsume all those grammatical relations (Subject-Predicate, Head-Modifier, etc.) under the concept of *dependency* relations.

Graph: a well understood model for representing objects that can be viewed as a connected set of more elementary sub-objects

Head: the most important word in a constituent; the word that carries the main meaning of the phrase. The head of a constituent cannot be left out.

Hierarchy: the relative position of constituents in a syntactic tree.

Human language technology: technology as applied to natural languages

Label: a feature value providing the content of an edge.

Main clause: a finite clause, which can act on its own as a complete sentence. Example: *The train has some delay*.

Modifier: a modifier is a part of the constituent which ascribes a property to the head of the constituent. A modifier may be placed before or after the head of the phrase (pre-modifier or post-modifier). Modifiers are optional in a constituent

Natural language processing NLP: field covering knowledge and techniques involved in the processing of linguistic data by a computer

Node: pair consisting of a (possibly multiple) span and a category,

Non-Terminal nodes have an outgoing constituency edges (to be defined)

Phrases: a word or group of words which can fulfil a *grammatical function* in a clause. But we allow empty phrases (as the example of the empty NP in Italian and Spanish, being non-referential pronouns and having the role of subjects in *clauses*). A phrase's typicality named after the most important word in it (which we also call the *head*), so we have for example noun phrases, verb phrases, adjectival phrases, adverbial phrases and prepositional phrases.). Phrases have been informally described as "bodied words", in that the parts of the phrase that are added to the head elaborate and specify the reference of the head word. In our model, a phrase is a special case of a constituent.

Sentences: a sequence of words, starting very often with a capital letter up to a final punctuation mark. But this definition is too restricted to layout property of certain language styles. A usage rule says that a complete sentence must contain a subject and a verb (in finite mode). A sentence consists of one or more *clauses*. In describing speech, it is common to talk about 'utterances' rather than sentences.

Span: a pair of points defining a segment of the document submitted to syntactic annotation. The first point is less or equal to the second point. A Multiple span is a sequence of spans where the ending point of each span is less or equal to the starting point of the subsequent span.

Specifier: a specifier in a constituent specifies the head (or the combination of modifier and head) with information about number, definiteness, proximity and ownership

Subordinated clause: a clause which fulfils a *grammatical function* in a *phrase* (for example a relative clause modifying the head noun of a nominal phrase) or in another clause. A subordinated clause can not act on its own as a *sentence*.

Syntax category on frame: set of restrictions and category properties of the syntactic arguments that can or must occur with it.

Example: A fred (syntactic argument) read a book (syntactic argument) today (adunct)

NOTE The subject, indirect object and direct object are possible grammatical categories for a sentence.

Syntax: The way words are grouped together in linguistic meaning units, and their relations that exist between those units.

Syntactic argument: one of the essential and functional elements in a clause that identifies the participants in the process referred to by a verb

Syntax Tree: a syntactic graph in which each node has a single parent.

Terminal node: refers to a single wordForm/excavation or a span with length =0, and the node and the wordForm/excavation have identical spans.

## 4 Key standards used by SynAF

### 4.1 Unicode

SynAF's Unicode component and presumes that all data are represented using Unicode character encodings.

### 4.2 ISO 12620 Data Category Registry (DCR)

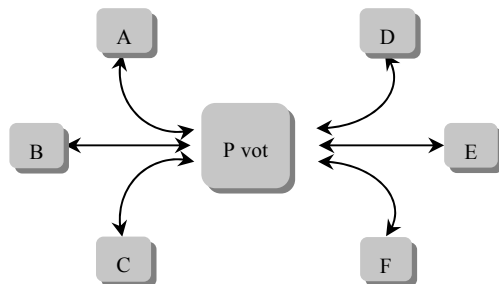
The designers of an SynAF conformant annotations use data categories from the ISO 12620 Data Category Registry (DCR), or a tagset that can be mapped onto the data categories.

### 4.3 Unified Modeling Language (UML)

SynAF complies with the specifications and modeling principles of UML as defined by the Object Management Group (OMG) [4]. SynAF uses a subset of UML that is relevant for linguistic description. (not done yet).

## 5 Embedding SynAF in the LAF model<sup>6</sup>

We want to embed the meta-model of SynAF in the more generic Linguistic Annotation Framework (LAF) and reuse its annotation strategy. LAF provides a general framework for representing annotations that has been described elsewhere in detail (Ide and Romary, 2004, 2006). Its development has been based on common practice and convergence of approaches in linguistic annotation over the past 15-20 years. The core of the framework is specification of an abstract model for annotations instantiated by a *pivot format*, into and out of which annotations are mapped for the purposes of exchange.



<sup>6</sup> The website section 5 is taken from (Ide, 2007).

**Figure 1: Use of the LAF pivot format**

Figure 1 shows the overall idea for six different user annotation formats (labelled A – F), which requires two mappings for each scheme—one into and one out of the pivot format, provided by the scheme designer. The maximum number of mappings among schemes is therefore  $2n$ , vs.  $n^2-n$  mutual mappings without the pivot.

To map to the pivot, an annotation scheme must be (or be rendered via the mapping) somorphic to the abstract mode, which consists of (1) a *referential structure* for associating stand-off annotations with primary data, instantiated as a directed graph; and (2) a *feature structure representation* for annotation content. An annotation thus forms a directed graph referencing  $n$ -dimensional regions of primary data as well as other annotations, in which nodes are labelled with feature structures providing the annotation content. Formally, LAF consists of:

- A data mode for annotations based on directed graphs defined as follows: A graph of annotations  $G$  is a set of vertices  $V(G)$ <sup>7</sup> and a set of edges  $E(G)$ . Vertices and edges may be labelled with one or more features. A feature consists of a quadruple  $(G', VE, K, V)$  where,  $G'$  is a graph,  $VE$  is a vertex or edge in  $G'$ ,  $K$  is the name of the feature and  $V$  is the feature value.
- A *base segmentation* of primary data that defines edges between virtual nodes located between each “character” in the primary data.<sup>8</sup> The resulting graph  $G$  is treated as an *edge graph*  $G'$  whose nodes are the edges of  $G$ , and which serve as the leaf (“s n”) nodes. These nodes provide the base for an annotation or several layers of annotation. Multiple segmentations can be defined over the primary data, and multiple annotations may refer to the same segmentation.
- Serializations of the data mode, one of which is designated as the pivot.
- Methods for manipulating the data mode.

Note that LAF does not provide specifications for annotation *content categories* (i.e., the labels describing the associated linguistic phenomena), for which standardization is a much trickier matter. The LAF architecture includes a *Data Category Registry* (DCR) containing pre-defined data elements and schemes that may be used directly in annotations, together with means to specify new categories and modify existing ones (see Ide and Romary, 2004).

## 6 The SynAF Metamodel

### 6.1 Introduction

While preparing SynAF, we identified some existing notations arising from some common data mode that seems to offer a good basis for the SynAF meta-mode (Tiger and ISST for example, but also a longer list of corpora has been studied, see Deverbeke D.3.1 of LIRICS). Based on this study we strongly suggest the adoption of a multi-layered annotation strategy interleaving syntactic annotation for both constituency and dependency in a sound representation scheme. The studied notations also offer a quite complete set of descriptors, which we started to “merge” into a first set of candidate data-categories, to be extended by data categories covering syntactic phenomena (constituency and dependency) for other languages than German and Italian. Our list of candidate data categories is presented in Annex A. TIGER and ISST are summarized in Annexes

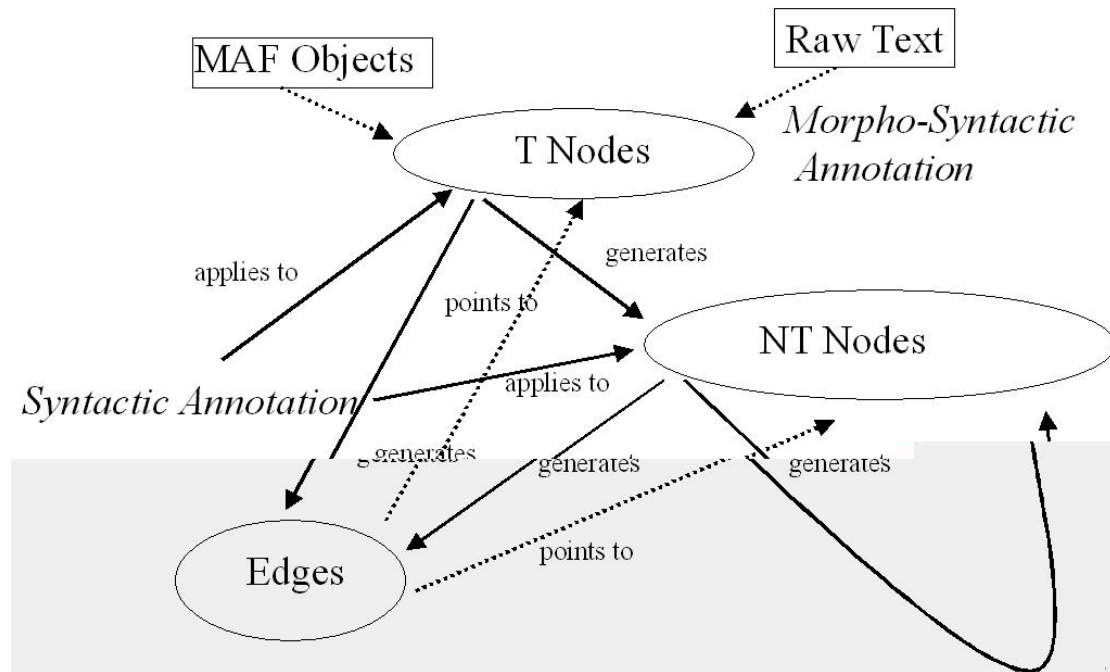
The SynAF mode will be represented by UML classes and by a set of ISO 12620 data categories that function as UML attribute-value pairs. The data categories are used to decorate

<sup>7</sup> The word “vertex” is a synonym to “node”.

<sup>8</sup> A character is defined to be a contiguous byte sequence of a specified length. For text, the default is UTF-16.

the UML classes that provide a graphical view of the model. SynAF specifications in the form of textual descriptions that describe the semantics of the modeling elements provide more complete information about the SynAF classes, relationships, and extensions that can be included in the UML diagram. Developers shall define a data category selection (DCS) as specified for SynAF data category selection procedures (see below).

**6.2 The SynAF diagram (to be represented in UML)**



**Figure 1: The SynAF metamodel**

**6.2.1 T Nodes class**

The *t\_nodes* class represents the terminal nodes of a syntax tree, mostly consisting of morpho-syntactically annotated words, but empty elements are allowed. The *t\_nodes* are defined over a *span*. They can be a multiple span (for accounting for discontiguous constituents). The *t\_nodes* are labeled with syntactic categories valid for the word level.

**6.2.2 NT Nodes class**

The *nt\_nodes* class represents the non-terminal nodes of a syntax tree, mostly consisting of *t\_nodes* and *nt\_nodes*, but empty elements are allowed. The *nt\_nodes* are also defined over a (possibly multiple) *span*. The *nt\_nodes* are labeled with syntactic categories valid at the phrase level and ger(c)ause, sentential).

**6.2.3 Edges class**

The *Edges* class represents the dependency relation between nodes (both terminal and non-terminal nodes). The dependency relation is a binary one and consists of a label name and a pair of source and target nodes.

**6.2.4 Syntactic Annotation class**

The *Syntactic Annotation* class represents the application of syntactic information to MAF annotated input. It can be either manual or an automatic application. When syntactic

## **ISO 24615:2007**

Annotations applied to nodes (non-terminal or terminal), then it generates either a new (non-terminal) node or a dependency edge.



## Annex A: (informative) Data Categories for SynAF

Our strategy consisted in collecting some of the most consensus syntactic annotation definitions for generating a list of data categories for constituency (node labels) and dependency (edge labels) annotation, which will be established in the document resulting from the work in ISO TC37/SC4 TDG 4 "Syntax". In this document we present the actual list of candidates, as they have been detected in annotation initiatives like TIGER, ISST, Sparce and EAGLES, and modified/ harmonized for the purpose of this document. We do not quote the specific origin of each candidate data category. We indicate, where appropriate, language specific data categories.

### A.1 Constituency

| Constituency_labels | Meaning  |
|---------------------|--|
| AA                  | superlative phrase with am (for German)                        |
| AP                  | adjective phrase   |
| AVP                 | adverbial phrase   |
| CAC                 | coordinated adposition   |
| CAP                 | coordinated adjective phrase                                   |
| CAVP                | Coordinated adverbial phrase                                   |
| CCP                 | Coordinated complementiser                                     |
| CH                  | Clause (non-recursive constituent)                             |
| CNP                 | Coordinated noun phrase  |
| CO                  | coordination   |
| CPP                 | Coordinated adpositional phrase                                |
| CVP                 | Coordinated verb phrase (non-finite)                           |
| CVZ                 | Coordinated infinitive with zu (for German)                    |
| NP                  | noun phrase  |
| PN                  | proper noun  |
| PP                  | adpositional phrase (prepositional and postpositional phrases) |
| S                   | Sentence   |
| VP                  | verb phrase (non-finite)                                       |
| VZ                  | infinitive with zu (for German)                                |

SPD                      prepositional phrase *di* 'of' (for Italian)

SPDA                    prepositional phrase *da* 'by, from' (for Italian)

IBAR                    verbal nucleus with finite tense and auxiliary

|       |                            |                   |
|-------|----------------------------|-------------------|
|       | elements                   | elements, adverbs |
|       | and negation               |                   |
| SV2   | infinitive clause          |                   |
| SV3   | participial clause         |                   |
| SV5   | gerundive clause           |                   |
| FAC   | sentential complement      |                   |
| FS    | subordinate sentence       |                   |
| FINT  | +wh interrogative sentence |                   |
| F2    | relative clause            |                   |
| CP    | dislocated or fronted      |                   |
|       | sentential adjuncts        |                   |
| COMPC | copulative/predicative     |                   |
|       | complement                 |                   |

## A.2 Dependency

In the following we present the candidate data categories for dependency structures (the abets of edges in the annotation graph). Source of inspiration here were the Spar e and the Tger tagsets for dependency. We use also some examples taken from Spar e (these sort below some data categories.)

**mod:** indicates the word introducing the dependent in a head-modifier relation

mod(of, gift, book)                      the gift of a book

mod(by, gift, Peter)                    the gift of a book by Peter

mod(of, examination, patient) the examination of the patient

mod('s, doctor, examination) the doctor's examination of the patient

**cmod, xmod, ncmo d:** Causal and non-causal modifiers may (optionally) be distinguished by the use of cmod / xmod, and ncmo d respectively, each with the same slots as **mod**. The GR cmo d s for when the adjunct is controlled from within, and xmo d for control from without the constituent under consideration.

cmod(because, eat, be)                ate the cake because he was hungry

xmod(without, eat, as)                ate the cake without asking

**subj:** indicates the subject in the grammatical relation Subject-Predicate. The relation between a predicate and its subject; where appropriate, the **initial\_gr** indicates the syntactic relation between the predicate and subject before any GR-changing process.

sub(arrive, John,\_)                    John arrived in Paris

sub(employ, Microsoft,\_) Microsoft employed 10 C programmers

sub(employ, Paul, ob)                Paul was employed by Microsoft

With pro-drop languages such as Italian, when the subject is not overtly realised the annotation s, for example, as follows:

sub(arrive, Pro,\_) arrive in ritardo '(I) arrived late'

Where the dependent is satisfied by the abstract filler **pro**, we can indicate that person and number of the subject can be recovered from the inflection of the head verb form.

**csubj, xsubj, ncsu bj:** The Grammatical Relations (RL) s **csubj** and **xsubj** may be used for causal subjects, controlled from within, or without, respectively. **ncsu bj** s a non-causal subject.

csub (leave, mean,\_) that Ne e left without saying good-bye meant she was still angry

## ISO 24615:2007

xsub (w n,require,\_) to w n t e Amer ca's Cup requ res eaps of cas

**dobj:** Indicates the object in the grammatical relation between a predicate and its direct object.

dob (read,boo ,\_) read boo s

dob (ma ,Mary, ob ) ma Mary t e contract

**iobj:** The relation between a predicate and a non-clausal complement introduced by a preposition; **type** indicates the preposition introducing the dependent.

ob ( n,arr ve,Spa n) arr ve n Spa n

ob ( nto,put,box) put t e too s nto t e box

ob (to,g ve,poor) g ve to t e poor

**obj2:** The relation between a predicate and the second non-clausal complement in ditransitive constructions.

ob 2( ead,dependent)

ob 2(g ve,present) g ve Mary a present

ob 2(ma ,contract) ma Pau t e contract

**dependent:** The most general relation between a head and a dependent

dependent( ntroducer, ead,dependent)

dependent( n, ve,Rome) Mar sa ves n Rome

| Dependency Rel               | ID  | Definition   | Parent               |
|------------------------------|-----|--|----------------------|
| Adpos tona Case<br>Mar er    | AC  | Prepos t on/postpos t on n a PP, annotated as a s ster<br>const tuent of t e determ ner, ad ect ves, noun etc        | PP                   |
| Ad ect ve<br>Component       | ADC | Component of a mu t-to en ad ect ve (MTA)  | MTA                  |
| Appos t on                   | APP | " nserted" p rase, furt er spec fy ng/mod fy ng t e<br>ent ty descr bed by t e matr x NP.                            | NP<br>PP             |
| Adverb a p rase<br>Component | AVC | Component of a ead- ess AVP  | ADV                  |
| con unct                     | CJ  | Const tuent part c pat ng n coord nat on   | any                  |
| comparat ve<br>con unct on   | CM  | L ngu st c part c es ntroduc ng a compar son n<br>comparat ve construct ons (for exam p e "grosser a s"<br>n German) |                      |
| dat ve                       | DA  | Dat ve ob ect/'free dat ve' (for nguages av ng t s<br>case n t e morp ogy/syntax))                                   | S<br>VP<br>AP<br>AVP |
| ead                          | HD  | T e ma n e ements n a nd of cons tuents  | S<br>VP<br>AP<br>AVP |
| postnom na mod f er          | MNR | Postnom na NP/PP mod f er  | NP<br>PP             |
| negat on                     | NG  | t e negat on part c e `n c t' (a so mod f ed)  | any                  |
| gen t ve ob ect              | OG  | Gen t ve ob ects of verbs, part c pes and certa n<br>ad ect ves (for nguage av ng t e gen t ve case n                |                      |

**ISO 24615:2007**

|                       |    |   |                           |
|-----------------------|----|---|---------------------------|
|                       |    | t e morp o ogy/syntax)  |                           |
| pred cate             | PD | Pred cat ve AP/NP/PP, typ ca y n a copu ar construct on                               | S<br>VP                   |
| morp o og ca part c e | PM | two cases: t e nf n t va `zu' (zu ge en) t e ad ect va (super at ve) `am' (am besten) | VZ AA                     |
| re at ve c ause       | RC |   | NP<br>PP<br>S<br>VP<br>AP |

## Annex B (informative) Annotation example

The following examples show how a multi-layered approach to syntactic annotation can be encoded in XML. The tagset in use is not pointing yet to the data categories, but such a thing will be included in the next version of the document. The Grammatical Functions (dependencies) are encoded here with the “cause” XML elements. The dependencies with constituents are not annotated explicitly in this example.

```
<?xml version='1.0' encoding='ISO-8859-1'?>
<document id="ww92-short.xml" lang="de">
<paragraph id="p1" corresp="">
  <sentence id="s4" corresp="">
    <clauses>
      <clause id="c11" from="c1" to="c22" pred_struct="c18+1"
pred_subcat_stem="sei" GF_Subj="c19+1" NP_List="c19+1"
VG_List="c18+1"/>
      <clause id="c12" from="c23" to="c26" pred_struct="c22+1"
pred_subcat_stem="find" GF_Subj="c24+1" GF_Acc_Obj="c23+1"
NP_List="c23+1 c24+1" VG_List="c22+1"/>
    </clauses>
    <chunks>
      <chunk id="c1" from="1" to="1" type="VG"/>
      <chunk id="c2" from="2" to="3" type="NP"/>
      <chunk id="c3" from="4" to="4" type="W"/>
      <chunk id="c4" from="5" to="5" type="W"/>
      <chunk id="c5" from="6" to="6" type="VG"/>
      <chunk id="c6" from="7" to="7" type="NP"/>
      <chunk id="c7" from="8" to="10" type="NP"/>
      <chunk id="c8" from="11" to="11" type="W"/>
    </chunks>
    <text>
      <token id="1" infl="[92 93 94 95]" pos="3" lemma="sei"
tc="22">Sind</token>
      <token id="2" infl="[2 5 20 6 13 23 9 16]" pos="8"
lemma="kein" tc="21">keine</token>
      <token id="3" infl="[6 7 8 9]" pos="1" lemma="angabe"
tc="22">Angaben</token>
      <token id="4" infl="[25]" pos="5" lemma="erhaeltlich"
tc="21">erhaeltlich</token>
      <token id="5" infl="[439]" pos="21" lemma=","
tc="1">,</token>
      <token id="6" infl="[204 205 209 206]" pos="2"
lemma="find" tc="21">findet</token>
      <token id="7" infl="[423 431 428 432 424 433 430 434]"
pos="11" lemma="sich" tc="21">sich</token>
      <token id="8" infl="[10 12]" pos="7" lemma="d-det"
tc="21">das</token>
      <token id="9" infl="[10 11 12 13 14 16]" pos="1"
lemma="kuerzel" tc="22">Kuerzel</token>
      <token id="10" tc="19">KA</token>
      <token id="11" infl="[440]" pos="21" lemma="."
tc="1">.</token>
```

```

</text>
</sentence>
<sentence id="s5" corresp="">
  <clauses>
    <clause id="c11" from="c1" to="c30" pred_struct="c27+1"
pred_subcat_stem="bedeut" GF_Subj="c26+1" NP_List="c26+1"
VG_List="c27+1"/>
  </clauses>
  <chunks>
    <chunk id="c1" from="1" to="2" type="NP"/>
    <chunk id="c2" from="3" to="3" type="VG"/>
    <chunk id="c3" from="4" to="14" type="SUBORD_CLAUSE"/>
    <chunk id="c4" from="15" to="15" type="W"/>
  </chunks>
  <text>
    <token id="1" infl="[17 10 12]" pos="8" lemma="ein"
tc="22">Ein</token>
    <token id="2" infl="[17 18 19]" pos="1" lemma="strich"
tc="22">Strich</token>
    <token id="3" infl="[204 205 209 107 206]" pos="2"
lemma="bedeut" tc="21">bedeutet</token>
    <token id="4" tc="1" lemma="," infl="[439]"
pos="21">,</token>
    <token id="6" tc="21" lemma="dass" pos="20">dass</token>
    <token id="8" tc="21" lemma="zutreff" infl="[204]"
pos="2">zutrifft</token>
    <token id="9" tc="21" lemma="nicht" pos="22">nicht</token>
    <token id="11" tc="22" lemma="wert" infl="[18 20 21 23]"
pos="1">Werte</token>
    <token id="12" tc="21" lemma="d-det" infl="[2 5 20 6 13 23
9 16]" pos="7">die</token>
    <token id="14" tc="21" lemma="nicht"
pos="22">nicht</token>
    <token id="16" tc="21" lemma="vergleichbar" infl="[25]"
pos="5">vergleichbar</token>
    <token id="18" tc="21" lemma="sei" infl="[92 93 94 95]"
pos="3">sind</token>
    <token id="20" tc="21" lemma="oder" pos="19">oder</token>
    <token id="22" tc="22" lemma="kriterium" infl="[10 11 12]"
pos="1">Kriterium</token>
    <token id="23" tc="21" lemma="d-det" infl="[10 12]"
pos="7">das</token>
    <token id="25" tc="22" lemma="spalt" infl="[18 20 21 23]"
pos="1">Spalte</token>
    <token id="26" tc="21" lemma="dies" infl="[17 3 4 21 7
14]" pos="7">dieser</token>
    <token id="28" tc="21" lemma="d-det" infl="[10 12]"
pos="7">das</token>
    <token id="29" tc="21" lemma="fuer" infl="[102]"
pos="23">fuer</token>
    <token id="30" tc="22" lemma="unternehmen" infl="[10 11 12
13 14 15 16]" pos="1">Unternehmen</token>
    <token id="31" infl="[440]" pos="21" lemma="."
tc="1">.</token>
  </text>
</sentence>

```

ISO 24615:2007

```

<sentence id="s6" corresp="">
  <clauses>
    <clause id="c11" from="c1" to="c37" pred_struct="c31+1"
    pred_subcat_stem="bezieh" GF_Subj="c30+1" GF_Acc_Obj="c32+1"
    PP_Adjunkt="c35+1" NP_List="c30+1 c32+1" PP_List="c33+1 c34+1
    c35+1" VG_List="c31+1"/>
  </clauses>
  <chunks>
    <chunk id="c1" from="1" to="3" type="NP"/>
    <chunk id="c2" from="4" to="4" type="VG"/>
    <chunk id="c3" from="5" to="5" type="NP"/>
    <chunk id="c4" from="6" to="8" type="PP"/>
    <chunk id="c5" from="9" to="10" type="PP"/>
    <chunk id="c6" from="11" to="11" type="PP"/>
    <chunk id="c7" from="12" to="12" type="W"/>
  </chunks>
  <text>
    <token id="1" infl="[2 5 20 6 13 23 9 16]" pos="7"
    lemma="d-det" tc="22">Die</token>
    <token id="2" infl="[41 42 43 44 45 46 47 48 49 50 51 52
    53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73
    74 75 76 77 78 79 80 81 82 83 84]" pos="6" lemma="angeb"
    tc="21">angegebenen</token>
    <token id="3" infl="[18 20 21 23]" pos="1" lemma="wert"
    tc="22">Werte</token>
    <token id="4" infl="[92 93 94 95 96 97 98 99 100 101]"
    pos="2" lemma="bezieh" tc="21">beziehen</token>
    <token id="5" infl="[423 431 428 432 424 433 430 434]"
    pos="11" lemma="sich" tc="21">sich</token>
    <token id="6" infl="[24 102]" pos="23" lemma="in"
    tc="21">in</token>
    <token id="7" infl="[17 3 4 21 7 14]" pos="7" lemma="d-
    det" tc="21">der</token>
    <token id="8" infl="[2 3 4 5]" pos="1" lemma="regel"
    tc="22">Regel</token>
    <token id="9" infl="[24 102]" pos="23" lemma="auf"
    tc="21">auf</token>
    <token id="10" infl="[20 21 23]" pos="1" lemma="abschluss"
    tc="22">Abschluesse</token>
    <token id="11" tc="11">31.</token>
    <token id="12" tc="11">1991</token>
    <token id="13" tc="21" lemma="zum" infl="[24]"
    pos="23">zum</token>
    <token id="14" tc="22" lemma="dezember" infl="[17 18 19 20
    21 22 23]" pos="1">Dezember</token>
    <token id="15" infl="[440]" pos="21" lemma="."
    tc="1">.</token>
  </text>
</sentence>
<sentence id="s7" corresp="">
  <clauses>
  </clauses>
  <chunks>
    <chunk id="c1" from="1" to="8" type="PP"/>
    <chunk id="c2" from="9" to="9" type="VG"/>
    <chunk id="c3" from="10" to="15" type="PP"/>
  </chunks>

```

## ISO 24615:2007

```
<chunk id="c4" from="16" to="16" type="VG"/>
<chunk id="c5" from="17" to="17" type="W"/>
</chunks>
<text>
  <token id="1" infl="[24 102]" pos="23" lemma="in"
tc="22">In</token>
  <token id="2" infl="[41 42 43 44 45 46 47 48 49 50 51 52
53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73
74 75 76 77 78 79 80 81 82 83 84]" pos="5" lemma="einig"
tc="21">einigen</token>
  <token id="3" infl="[22]" pos="1" lemma="fall"
tc="22">Faellen</token>
  <token id="4" pos="20" lemma="wie" tc="21">wie</token>
  <token id="5" pos="22" lemma="beispielsweise"
tc="21">beispielsweise</token>
  <token id="6" tc="22" lemma="siemens" infl="[10 0 11 12 13
14 16]" pos="1">Siemens</token>
  <token id="7" pos="19" lemma="oder" tc="21">oder</token>
  <token id="8" tc="19">MAN</token>
  <token id="9" infl="[204]" pos="2" lemma="werd"
tc="21">wird</token>
  <token id="10" infl="[102]" pos="23" lemma="per"
tc="21">per</token>
  <token id="11" infl="[10 11 12]" pos="1" lemma="ende"
tc="22">Ende</token>
  <token id="12" infl="[17 18 19 20 21 23]" pos="1"
lemma="september" tc="22">September</token>
  <token id="13" pos="17" lemma="beziehungsweise"
tc="21">beziehungsweise</token>
  <token id="14" infl="[10 11 12]" pos="1" lemma="ende"
tc="22">Ende</token>
  <token id="15" infl="[17 18 19]" pos="1" lemma="juni"
tc="22">Juni</token>
  <token id="16" infl="[204 205 107 206]" pos="2"
lemma="bilanzier" tc="21">bilanziert</token>
  <token id="17" infl="[440]" pos="21" lemma="."
tc="1">.</token>
</text>
</sentence>
...
</Paragraph>
</document>
```



## Bibliography

- Abe é A., S. Hansen-Sc rra & H. Usz ore t (eds.), 2003. Proceed ngs of t e 4t Internat ona Wor s op on L ngu st ca y Interpreted Corpora (LINC-03).
- Ca zo ar N., Mc Naug t J., Zampo A. 1996 Eag es, ed tors ntroduct on. <http://www.cnr.it/EAGLES96/edntro/edntro.htm>
- Ca zo ar N., Bertagna F., Lenc A., Monac n M. ed tors, 2003. Standards and best Pract ce for Mu t ngua Computat ona Lex cons. MILE (T e Mu t ngua ISLE Lex ca Entry). ISLE CLWG De verab e D2.2 & 3.2 P sa.
- Rumbaug J., Jacobson I., Booc G. T e un fed mode ng language reference manua , second ed t on Addison Wesley 2004
- S. Montemagn , F. Barsott , M. Batt sta, N. Ca zo ar , A. Lenc O. Corazzar , A. Zampo , F. Fanc u , M. Massetan , R. Bas R. Raffae , M.T. Paz enza, D. Sarac no, F. Zanzotto, F. P anes N. Mana, and R. De monte. Bu d ng t e Ita an Syntact c-Semant c Treeban . In Anne Abe é, ed tor, Bu d ng and Us ng syntact ca y annotated corpora, pages 189--210. K uwer, Dordrec t, 2003.
- Nancy Ide and Laurent Romary, 2004. A Registry of Standard Data Categories for Linguistic Annotation. Proceed ngs of t e Fourt Language Resources and Eva uat on Conference (LREC), Lisbon, 135-39.
- Nancy Ide and Laurent Romary, 2004. Internat ona Standard for a Linguistic Annotation Framework . Journa of Natura Language Eng neer ng, 10:3-4, 211-225.
- Nancy Ide and Laurent Romary, 2006. Represent ng Linguistic Corpora and T er Annotations. Proceed ngs of t e F ft Language Resources and Eva uat on Conference (LREC), Genoa, Ita y.
- Nancy Ide, GrAF: A Grap -based Format for Linguistic Annotations. Proceed ngs of t e LAW Wor s op at ACL 2007, Prague.
- T e EAGLES nt at ve: <http://www.cnr.it/EAGLES96/home.htm>
- T e LIRICS Pro ect: <http://r.cs.or.a.fr>
- T e SPARKLE Pro ect: <http://www.cnr.it/spare/spare.htm>
- T e TIGER pro ect: <http://www.ms.un-stuttgart.de/projekte/TIGER/TIGERCorpus/>