

SVENSK STANDARD

SS-ISO 24611:2015



Fastställt/Approved: 2015-03-16
Publicerad/Published: 2015-03-18
Utgåva/Edition: 1
Språk/Language: engelska/English
ICS: 01.020

Hantering av språkliga resurser – Morfosyntaktiskt Annotationsformat (MAF) (ISO 24611:2012, IDT)

Language resource management – Morpho-syntactic annotation framework (MAF) (ISO 24611:2012, IDT)

This preview is downloaded from www.sis.se. Buy the entire standard via <https://www.sis.se/std-8013206>

Standarder får världen att fungera

SIS (Swedish Standards Institute) är en fristående ideell förening med medlemmar från både privat och offentlig sektor. Vi är en del av det europeiska och globala nätverk som utarbetar internationella standarder. Standarder är dokumenterad kunskap utvecklad av framstående aktörer inom industri, näringsliv och samhälle och befrämjar handel över gränser, bidrar till att processer och produkter blir säkrare samt effektiviserar din verksamhet.

Delta och påverka

Som medlem i SIS har du möjlighet att påverka framtida standarder inom ditt område på nationell, europeisk och global nivå. Du får samtidigt tillgång till tidig information om utvecklingen inom din bransch.

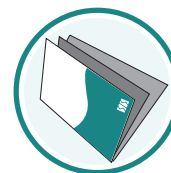
Ta del av det färdiga arbetet

Vi erbjuder våra kunder allt som rör standarder och deras tillämpning. Hos oss kan du köpa alla publikationer du behöver – allt från enskilda standarder, tekniska rapporter och standardpaket till handböcker och onlinetjänster. Genom vår webbtjänst e-nav får du tillgång till ett lättnavigerat bibliotek där alla standarder som är aktuella för ditt företag finns tillgängliga. Standarder och handböcker är källor till kunskap. Vi säljer dem.

Utveckla din kompetens och lyckas bättre i ditt arbete

Hos SIS kan du gå öppna eller företagsinterna utbildningar kring innehåll och tillämpning av standarder. Genom vår närhet till den internationella utvecklingen och ISO får du rätt kunskap i rätt tid, direkt från källan. Med vår kunskap om standarders möjligheter hjälper vi våra kunder att skapa verklig nytta och lönsamhet i sina verksamheter.

Vill du veta mer om SIS eller hur standarder kan effektivisera din verksamhet är du välkommen in på www.sis.se eller ta kontakt med oss på tel 08-555 523 00.



Standards make the world go round

SIS (Swedish Standards Institute) is an independent non-profit organisation with members from both the private and public sectors. We are part of the European and global network that draws up international standards. Standards consist of documented knowledge developed by prominent actors within the industry, business world and society. They promote cross-border trade, they help to make processes and products safer and they streamline your organisation.

Take part and have influence

As a member of SIS you will have the possibility to participate in standardization activities on national, European and global level. The membership in SIS will give you the opportunity to influence future standards and gain access to early stage information about developments within your field.

Get to know the finished work

We offer our customers everything in connection with standards and their application. You can purchase all the publications you need from us - everything from individual standards, technical reports and standard packages through to manuals and online services. Our web service e-nav gives you access to an easy-to-navigate library where all standards that are relevant to your company are available. Standards and manuals are sources of knowledge. We sell them.

Increase understanding and improve perception

With SIS you can undergo either shared or in-house training in the content and application of standards. Thanks to our proximity to international development and ISO you receive the right knowledge at the right time, direct from the source. With our knowledge about the potential of standards, we assist our customers in creating tangible benefit and profitability in their organisations.

If you want to know more about SIS, or how standards can streamline your organisation, please visit www.sis.se or contact us on phone +46 (0)8-555 523 00



Den internationella standarden ISO 24611:2012 gäller som svensk standard. Detta dokument innehåller den officiella engelska versionen av ISO 24611:2012.

The International Standard ISO 24611:2012 has the status of a Swedish Standard. This document contains the official English version of ISO 24611:2012.

© Copyright/Upphovsrätten till denna produkt tillhör SIS, Swedish Standards Institute, Stockholm, Sverige. Användningen av denna produkt regleras av slutanvändarlicensen som återfinns i denna produkt, se standardens sista sidor.

© Copyright SIS, Swedish Standards Institute, Stockholm, Sweden. All rights reserved. The use of this product is governed by the end-user licence for this product. You will find the licence in the end of this document.

Uppllysningar om sakinnehållet i standarden lämnas av SIS, Swedish Standards Institute, telefon 08-555 520 00. Standarder kan beställas hos SIS Förlag AB som även lämnar allmänna upplysningar om svensk och utländsk standard.

Information about the content of the standard is available from the Swedish Standards Institute (SIS), telephone +46 8 555 520 00. Standards may be ordered from SIS Förlag AB, who can also provide general information about Swedish and foreign standards.

Denna standard är framtagen av kommittén för Terminologi och språkliga resurser, SIS/TK 115.

Har du synpunkter på innehållet i den här standarden, vill du delta i ett kommande revideringsarbete eller vara med och ta fram andra standarder inom området? Gå in på www.sis.se - där hittar du mer information.

Contents

Page

Foreword	v
Introduction.....	vi
1 Scope	1
2 Normative references	1
3 Terms and definitions	1
4 The MAF meta-model	4
4.1 Overview.....	4
4.2 MAF Meta-model.....	4
5 Segmenting with tokens	6
5.1 General	6
5.2 Formal description: <token>	7
5.3 Embedding notation.....	7
5.4 Alternate representation for TEI based documents.....	8
5.5 Stand-off notation.....	9
5.6 Informative attributes.....	9
5.7 Completing the inline token notation	10
5.7.1 Joining tokens in embedded mode	10
5.7.2 Overlapping tokens	11
6 Word-forms as linguistic units.....	11
6.1 Formal description: <wordForm>	12
6.2 Token attachment.....	12
6.2.1 One token; one word-form	12
6.2.2 Several contiguous tokens; one word-form	12
6.2.3 Several discontinuous tokens; one word-form.....	13
6.2.4 Zero token; one word-form.....	13
6.2.5 One token; several word-forms	14
6.3 Referring to lexical entries	14
6.4 Compound word-forms.....	15
6.5 Identification of word-forms within a TEI-compliant document	15
7 Morpho-syntactic content.....	18
7.1 General	18
7.2 Using feature structures	18
7.3 Compact morpho-syntactic tags	18
7.4 FSR libraries	19
7.5 Designing tagsets.....	20
7.6 Formal description: <tagset>	22
8 Handling ambiguities	22
8.1 Word-form content ambiguities	22
8.2 Lexical Ambiguities.....	23
8.3 Structural ambiguities.....	23
8.3.1 Structural ambiguities with word-forms	23
8.3.2 Structural ambiguities with tokens.....	24
8.4 Simplified structuring variants	24
8.4.1 Non-ambiguous linear representation	24
8.4.2 Mixed linear and lattice representation.....	25
8.5 Expanding the simplified variants	26
8.5.1 Separating tokens and word-forms	26
8.5.2 Wrapping into local lattices.....	26

8.5.3	Merging local lattices	27
8.5.4	Removing <wfAlt>.....	28
8.6	Formal description: <wfAlt> and <fsm>	29
Annex A (informative) Encoded example using the MAF serialization.....		30
Annex B (normative) MAF specification		33
B.1	Elements	33
B.1.1	<dcs/>.....	33
B.1.2	<fsm>	34
B.1.3	<maf>	34
B.1.4	<tagset>	35
B.1.5	<token>	35
B.1.6	<transition>	36
B.1.7	<wfAlt>	36
B.1.8	<wordForm>	37
B.2	Model classes.....	38
B.3	Attribute classes	38
B.3.1	att.token.information	38
B.3.2	att.token.join.....	39
B.3.3	att.token.span.....	39
B.3.4	att.wordForm.content.....	39
B.3.5	att.wordForm.tokens	40
B.4	Macros	40
B.4.1	data.certainty.....	40
B.4.2	data.code	40
B.4.3	data.count.....	40
B.4.4	data.duration.w3c	41
B.4.5	data.enumerated	41
B.4.6	data.key.....	41
B.4.7	data.language.....	42
B.4.8	data.name	43
B.4.9	data.numeric.....	43
B.4.10	data.pointer	43
B.4.11	data.probability	44
B.4.12	data.temporal.w3c.....	44
B.4.13	data.truthValue.....	44
B.4.14	data.word	45
B.4.15	data.xTruthValue.....	45
Annex C (normative) Morpho-syntactic data categories		46
Bibliography		58

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, in liaison with ISO, also take part in the work. 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 2.

The main task of technical committees is to prepare International Standards. 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.

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO shall not be held responsible for identifying any or all such patent rights.

ISO 24611 was prepared by Technical Committee ISO/TC 37, *Terminology and other language and content resources*, Subcommittee SC 4, *Language resource management*.

Introduction

ISO/TC 37/SC 4 focuses on the definition of models and formats for the representation of annotated language resources. To this end, it has generalised the modelling strategy initiated by its sister committee, SC 3, for the representation of terminological data [Romary, 2001], through which linguistic data models are seen as the combination of a generic data pattern (a meta-model), which is further refined through a selection of data categories that provide the descriptors for this specific annotation level. Such models are defined independently of any specific formats, and ensure that an implementer has the necessary conceptual instrument with which to design and compare formats with regard to their degrees of interoperability.

One important aspect of representing any kind of annotation is the capacity to provide a clear and reliable semantics for the various descriptors used, either in the form of formal features and feature values, or directly as objects in a representation that is expressed, for instance, in XML. In order to be shared across various annotation schemas and encoding applications, such a semantics should be implemented as a centralised registry of concepts: we will henceforth refer to these as data categories. As such, data categories should bear the following constraints.

- From a technical point of view, they must provide unique, stable references (implemented as persistent identifiers, in the sense of ISO 24619) such that the designer of a specific encoding schema can refer to them in his or her specification. By doing so, two annotations will be deemed to be equivalent when they are in fact defined in relation to the same data categories (as feature and feature value).
- From a descriptive point of view, each unique semantic reference should be associated with precise documentation combining a full text elicitation of the meaning of the descriptor with the expression of specific constraints that bear upon the category.

In recent years, ISO has developed a general framework for representing and maintaining such a registry of data categories, encompassing all domains of language resources. This initiative, described in ISO 12620, has led to the implementation of an online environment providing access to all data categories that have been standardized in the context of the various language resource-related activities within ISO, or specifically as part of the maintenance of the data category registry. It also provides access to the various data categories that individual language technology practitioners have defined in the course of their own work and decided to share with the community.

The ISO data category registry, as available through the ISOCat (www.isocat.org) implementation, is intended as a 'flat' marketplace of semantic objects, providing only a limited set of ontological constraints. The objective there is to facilitate the maintenance of a comprehensive descriptive environment where new categories are easily inserted and reused without the need for any strong consistency check with the registry at large. Indeed, the following basic constraints are part of the data category model, as defined in ISO 12620:

- simple generic-specific relations, when these are useful for the proper identification of interoperability descriptors between data categories. For instance, the fact that `/properNoun/` is a sub-category of `/noun/` makes it possible to compare morpho-syntactic annotations based on different descriptive levels of granularity;
- the description of conceptual domains, in the sense of ISO 11179, to identify, when known or applicable, the possible value of so-called complex data categories. For instance, it can be used to record that possible values of `/grammaticalGender/` (limited to a small group of languages [Romary 2011]), could be a subset of `{/masculine/, /feminine/ and /neutral/}`;
- language-specific constraints, either in the form of specific application notes or as explicit restrictions bearing upon the conceptual domains of complex data categories. For instance, it is possible to express explicitly that `/grammaticalGender/` in French can only take the two values: `{/masculine/ and /feminine/}`.

This International Standard provides a comprehensive framework for the representation of morpho-syntactic (also referred to as part-of-speech) annotations. Such an annotation level corresponds to a first lexical abstraction level over language data (textual or spoken) and, depending on the language to be annotated, together with the characteristics of the annotation tool or annotation scheme that is being used, can vary enormously in structure and complexity.

In order to deal with such complex issues as ambiguity and determinism in morpho-syntactic annotation, this International Standard introduces a meta-model that draws a clear distinction between the two levels of tokens (representing the surface segmentation of the source) and word-forms (identifying lexical abstractions associated with groups of tokens). These two levels share the following specificities: on the one hand, they can be represented as simple sequences and as local graphs such as multiple segmentations and ambiguous compounds; on the other hand, any n-to-n combination can stand between word forms and tokens.

As linguistic segments (sometimes called ‘markables’ in the literature [see, for instance, Carletta et al. 1997]), *tokens* may be embedded in the source document as inline mark-up, or they may point remotely to it by means of so-called stand-off annotations.

As linguistic abstractions, *word-forms* can be qualified by various linguistic features characterising the morpho-syntactic properties that are instantiated in the realisation of the lexical entry within the annotated text. Such properties may range from the simple indication of a lemma up to an explicit reference to a lexical entry in a dictionary. In most existing applications of morpho-syntactic annotation, linguistic properties are expressed by means of so-called tags; these codes refer to basic feature structures (see early examples in Monachini and Calzolari, 1994). Such codes may also provide morphological information, including its part of speech (e.g. noun, adjective or verb), and features such as number, gender, person, mood and verbal tense.

In keeping with the general modelling strategy of ISO/TC 37, this International Standard/MAF provides means of relating morpho-syntactic tags expressed as feature structures (compliant with ISO 24610) to the data categories available in ISOCat. A normative annex of this International Standard elicits a core set of data categories that can be used as reference for most current morpho-syntactic annotation tasks in a multilingual context. However, when implementers of this International Standard find these categories inappropriate in either coverage, scope or semantics, they are encouraged to use ISOCat to define their own categories in compliance with ISO/TC 37 principles.

Associated to the meta-model, MAF also provides a default XML syntax that may be used to serialise MAF-compliant annotation models. Since many existing projects are based on the text encoding initiative (TEI) guidelines (www.tei-c.org) — particularly in digital humanities, where a proper encoding of textual sources is essential — this International Standard will also provide clues about how to articulate the MAF model with TEI-compliant encodings. Indeed, the TEI guidelines already offer a variety of constructs and mechanisms to cope with many issues relevant to spoken corpora and their annotations (Romary and Witt, 2012).

Finally, it should be noted here that this International Standard forms the conceptual basis for the development of the ISO 24614 series on word segmentation, whereby all general principles and rules defined in ISO 24614-1, as well as the constraints expressed in additional parts for specific languages, are to be understood according to the token–word-form dichotomy.

Language resource management — Morpho-syntactic annotation framework (MAF)

1 Scope

This International Standard provides a framework for the representation of annotations of word-forms in texts; such annotations concern tokens, their relationship with lexical units, and their morpho-syntactic properties.

It describes a metamodel for morpho-syntactic annotation that relates to a reference to the data categories contained in the ISOCat data category registry (DCR, as defined in ISO 12620). It also describes an XML serialization for morpho-syntactic annotations, with equivalences to the guidelines of the TEI (text encoding initiative).

2 Normative references

The following referenced documents are indispensable for the application of this document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO 24610-1, *Language resource management — Feature structures — Part 1: Feature structure representation*

3 Terms and definitions

For the purposes of this document, the terms and definitions given in ISO 24610-1 and the following apply.

3.1

DAG

directed acyclic graph

graph with directed edges and no cycles

Note 1 to entry: DAGs are a subset of *finite state automata* (3.4).

3.3

feature structure

set of feature specifications, used in the morpho-syntactic annotation framework (MAF) to express morpho-syntactic content

Note 1 to entry: Feature structures are described in ISO 24610-1.

3.4

FSA

finite state automata

graphs made up of states with an initial state and a final state, and a finite set of transitions from state to state

Note 1 to entry: See also *DAG* (3.1).