## ISO TC 37/SC 4 SC4N119.pdf Rev. 2
# Language Resource Management

Descriptors and Mechanisms for Language Resources

| | |
|---|---|
| Title: | Draft - Language Resource Management – Morpho-Syntactic Annotation Framework (MAF) |
| Editor(s): | Lionel Clément & Eric de la Clergerie |
| Source: | WG2 |
| Project number: | NP 24611 |
| Status: | WD |
| Date: | July 1, 2004 |
| Agenda/Action: | For Review |
| References: | |

# Language Resource Management –
# Morpho-Syntactic Annotation Framework (MAF)

Language Resource Management –
Morpho-Syntactic Annotation Framework (MAF)

# Introduction

TO BE COMPLETED

# Scope

In Natural Language Processing (NLP), the morpho-syntactic annotation phase assigns to each document segment (either text or speech) one or more *tags* providing information about the *part of speech* (noun, adjective, verb, . . . ), morphological and grammatical features (such as number, gender, person, mode, verbal tense, . . . ) and possibly other specific linguistic properties. The morpho-syntactic annotations attached to a segment do not refer to other segments or annotations, even if the choice of an annotation may depend on the surrounding context.

# Normative references

- Data Category Registry (DCR)

- Feature Structure Representation (FSR) and Feature Structure Declaration (FSD)

- Linguistic Annotation Framework (LAF)

- Text Encoding Initiative (TEI) – Chapters to be precised

- MPEG7 – About referring positions in multimedia documents

# Terms and definitions

### Associative relation

Relations by which a linguistic unit is associated with others. This is a mental association which does not requires their effective presence.[1]

### DAG - Directed Acyclic Graph TO BE COMPLETED

### FSA - Finite State Automata TO BE COMPLETED See DAG.

### Lemma

Class of inflected forms differing only by inflectional morphology. A lemma is usually refered to by one of these forms, arbitrarily chosen (e.g. infinitive for french verbs).

### Lexeme

Lexical **morpheme**. Distinguished from a grammatical morpheme by the fact that it belongs to an open list and that it bears an autonomous signification.

---

[1]It differs from a paradigmatic relation because the latter only refers to linguistic units associated by substitutability.

**Morpheme**

Smallest linguistic entity bearing a signification in a discourse. A morpheme is either grammatical (grammeme) or lexical (lexeme).

**Natural Language Processing**

Field covering knowledge and techniques allowing computerized processing of linguistic data.

This field requests skills related, among others, to linguistics, mathematical logic, statistics, and algorithmics.

**Syntagmatic relation**

Relations by which linguistic units present in the discourse are associated.

**Morpho-syntactic tag**

To an associative relation corresponds a feature, for which the related entities share the same value. The morpho-syntactic tag lists some of these features (part-of-speech, grammatical category, etc.).

**Token**

Connex and not-empty discourse sequence identified as such by a morpho-phonologic analysis or an automatic processing of the discourse.

This can involve the recognition of a regular or algebric language (matching of the separators), or a lexicological analysis (recognition of roots, morphological derivation and inflection, etc.).

**Morpho-syntactic unit - Word-form**

Connex or not-connex entity from a speech sequence identified as such in an **associative relation**. This identification is the basis of morpho-syntactic tagging (part-of-speech, grammatical category, agreement feature, etc.). Morpho-syntactic units may have no acoustic or graphic realization, or correspond to one or more **tokens**.

**Word Lattice** TO BE COMPLETED See DAG.

# Contents

# 1   General characteristics

## 1.1   Overview

In the Linguistic community, Morpho-Syntactic Annotations provide an important layer of linguistic information to a document, even if they do not cover the full range of possible linguistic annotations. For instance, other kind of annotations on references, discourse, prosody, or parsing may complete morpho-syntactic annotations.

Syntax and semantic can not be avoided in the definition of the parts of speech and of the grammatical categories. For instance, pronouns and substantives intrinsically carry a reference (to some entity); the tense or the aspect of verbs indicate the temporal deixis; the person, modality and other grammatical categories indicate the enunciation situation, . . . .

Therefore, it is not easy to provide an exact and precise definition of what cover morpho-syntactic annotations because they are strongly related to many other linguistic properties for a given language and for a given context.

Nevertheless, the present proposal tries to delimit minimal and maximal sequences in documents (either text or speech) that can be identified as *morpho-syntactic units* and tries to categorize the linguistic properties that may be used to mark these units, within some larger syntagmatic context. Minimal units can not be broken into subparts that could be identified by similar morpho-syntactic criteria, but may however still be broken into more atomic subparts with morphological or phonological properties. Morpho-syntactic units can be nested to form maximal units (such as compound words) that act as elementary units for other level of linguistic analysis, in particular parsing. The exact boundary between morpho-syntax and parsing may sometimes be fuzzy.
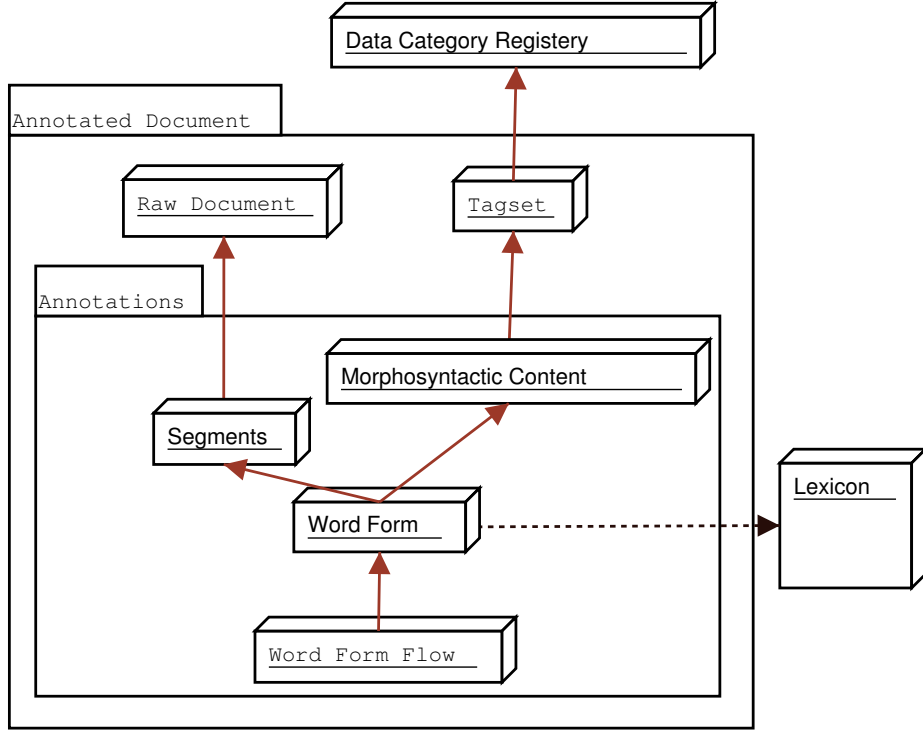
## 1.2 Meta-Model



Figure 1: Simplified view of MSAF meta-model

Figure 1 presents a simplified view of the proposed meta-model for morpho-syntactic annotations. An annotated document is formed by a raw document and a set of annotations. The annotations are carried by *word forms* covering zero, one or more segments or `tokens` of the documents. The word forms are organized in one or more flows, because of ambiguities. These flows are materialized by *finite state automata*. A word form may reference a lexicon entry. The morpho-syntactic content attached to a word form is expressed by *feature structures* following the guidelines of one or more *tagsets*. The terminology or set of *categories* used in tagsets are described w.r.t. *registered data categories*. The current proposal addresses the representation of segments (through tokens), word forms, word form flows, morpho-syntactic content, and tagsets.

## 2 Structuring

Morpho-syntactic annotations preserve the linear characteristic of the «signifiant» (i.e. the annotated document), either in space or time, depending on the nature of the annotated document. These documents correspond to linear sequences, be they texts, vocal records, transcriptions, ....

Therefore, the sequences that may be decorated by a morpho-syntactic annotation are not necessarily characterized by character strings but more generally by some *span* (or *interval*) in various kinds of documents. This dissociation of the morpho-syntactic annotations from the original document may be also useful when these annotations conflict with some

already present structuring of the document or when other kinds of annotation are to be attached that do not follow the same segmentation.

The *Text Encoding Initiative* (TEI) provides various mechanisms to identify sequences within documents and to link these sequences, as described in section 14.3 (*Blocks, Segments and Anchors*).

The current proposal for normalizing morpho-syntactic annotations should complete these mechanisms. In particular, it tries not to raise contradictions with TEI, in order to be able to mix TEI annotations and morpho-syntactic annotations. More generally, the possibility to easily add other levels of annotations (with their own segmentation) should be taken into account.

For instance, the segmentation of a document addresses positions in the original document that may be modified (insertion, displacement, deletion, . . . ). *TEI (14.1.1 Pointers and Links)* provides such a mechanism used to identify hypertextual links and other cross references with a document:

> `<ptr>` defines a pointer to another location in the current document in terms of one or more identifiable elements. target specifies the destination of the pointer by supplying the values used on the id attribute of one or more other elements in the current document

In order to reuse existing standards, such a mechanism will be used when possible to delimit linear sequences within textual documents. However, other kinds of position addressing will be considered for other kinds of documents (such a vocal recording), based on other existing standards (such as *MPEG7*).

## 2.1 Segmenting with tokens

Morpho-syntactic annotations are carried by segments *in praesentia* in the document flow, but this does not imply that the resulting segmentation corresponds to a sequence of adjacent segments partitioning the original document. It is particularly important to distinguish the morpho-syntactic units from their realizations. Some parts of a document may carry no annotations (typographic marks, didascalies, markup elements, . . . ); other parts may not exactly correspond to their segmented form (abbreviations, brachygraphies, typographic errors and variations, typographic and morphological contractions, . . . ). Also, a morpho-syntactic unit may not correspond to a segment identified by typographic marks (such as white spaces or hyphens), for instance for German compound words, speech transcription, or Sanskrit writing.

The element `token` is used to represent these segments of the original document that, roughly speaking, follow typographical, morphological, or phonological boundaries. The current proposal does not precise the linguistic properties of tokens. Among languages, a token may be identified through typographic properties (white-space, hyphens, characters, . . . ) and/or morphological properties (radical, affix, morpheme, . . . ). In any case, the description of the morphological, phonological or lexicologic structures that may define a token is not covered by the current proposal.

Other typographical marks used to format pages or to separate words and paragraphs, as well as encoding information, do not belong to morpho-syntactic annotations and are also not covered by this proposal, but rather by TEI.

In other words and as already mentioned, the element `token` provides an independence from the original document w.r.t. these aspects by providing a way to reference intervals in documents. The use of the attributes `from` and `to` are used to define such intervals. The content of these attributes depends on some chosen addressing schema. The following example shows an addressing schema based on the use of TEI element `ptr` (TEI Section 14.1.1 «Pointers and Links»):

```
  <s><ptr target="p1"/>The <ptr target="p2"/>victim
2 <ptr target="p3"/>'<ptr target="p4"/>s
  <ptr target="p5"/>friends <ptr target="p6"/>told
4 <ptr target="p7"/>police <ptr target="p8"/>that
  <ptr target="p9"/>Krueger <ptr target="p10"/>drove
6 <ptr target="p11"/>into <ptr target="p12"/>the
  <ptr target="p13"/>quarry <ptr target="p14"/>and
8 <ptr target="p15"/>never <ptr target="p16"/>surfaced
  <ptr target="p17"/>.<ptr target="p18"/></s>
```

```
  <token id="t1" from="p1" to="p2"/>
2 <token id="t2" from="p2" to="p4"/>
  <token id="t3" from="p4" to="p5"/>
4 <token id="t4" from="p5" to="p6"/>
  ...
```

It is not always necessary to separate the original document from its annotations. For simple cases, when one wishes to directly enrich a document with morpho-syntactic annotations, an alternate notation for `token` is provided that can embed textual content.

```
  <token id="t1">The</token>
2 <token id="t2">victim</token>
  <token id="t3">'s</token>
4 <token id="t4">friends</token>
  <token id="t5">told</token>
6 <token id="t6">police</token>
  <token id="t7">that</token>
8 <token id="t8">Krueger</token>
  <token id="t9">drove</token>
10 <token id="t10">into</token>
  <token id="t11">the</token>
12 <token id="t12">quarry</token>
  <token id="t13">and</token>
14 <token id="t14">never</token>
  <token id="t15">surfaced</token>
16 <token id="t16">.</token>
```

This alternate notation can only be used when the morpho-syntactic annotations do no conflict with other kinds of annotations. Furthermore, the content of the textual material separating the textual content embedded within `token` is not precisely defined (white-space, newlines, no space, hyphen, ...).

**QUESTION:** Should we add some way to precise this inter token material, for instance with an attribute `glue` on `token`.

Two tokens may overlap on some part of the original document, for instance to denote some agglutinated or contracted form (such as, in French, «des» as a contraction for «de

les» [of the]), or to denote multi-locutor documents with overlapping. In these cases, a **token** does not mark the realization of a typographical or vocal sequence, but expresses a deeper linguistic reality pertinent for segmenting a document.

Tokens address segments of the original document but also provide a level of possible abstraction w.r.t. this document, for instance w.r.t. graphical or phonological variations that are not linguistically pertinent. The non mandatory attribute **value** can be used to perform this abstraction, providing, for instance, the phonetic transcription of a speech segment, the roman transliteration of some Cyrillic word, the extension of an abbreviation, the corrected form of a words with typos, or the choice of a normalized form in presence of variations:

```
  <token value="et␣caetera" id="t1">etc.</token>
2 <token value="tzar" id="t2">csar</token>
  <token value="tzar" id="t3">tsar</token>
4 <token value="23/02/03" id="t4">February, 23rd 2003</token>
```

**META: Add an example of phonetic transcription**

The attribute **value** holds the linguistic interpretation of the content of a segment, pertinent to anchor morpho-syntactic annotations. Of course, this interpretation may differ from an user (or tool) to another one.

The abstraction provided by the attribute **value** is also adequate to handle the phenomena of contraction and agglutination where two tokens may cover the same segment of the original document for distinct values.

For instance, the following example illustrates the contraction of an abbreviation with a punctuation mark for «etc.», for the two variants of notations for element **token**:

```
a
    <token value="et␣caetera" id="t1" from="p1" to="p3"/>
2   <token value="#dot#" id="t2" from="p1" to="p3"/>
```

```
b
    <token value="et␣caetera" id="t1">etc.</token>
2   <token value="#dot#" id="t2"/>
```

**META: Should we consider a reference notation such as**

```
    <token value="et␣caetera" id="t1">etc.</token>
2   <token value="#dot#" id="t2" idref="t1"/>
```

Another example, in Modern Greek, is provided by the idiomatic expression "καλόκαγαθος" (*good and brave*) that may be segmented in three agglutinated segments "καλός", "και", and "αγαθος" and represented by:

```
  <token value="καλός" id="t0">καλο</token>
2 <token value="και" id="t1">κ</token>
  <token value="αγαθός" id="t2">αγαθος</token>
```

**QUESTION:** The representation of overlapping tokens with the embedding notation is slightly confusing (see "etc." example), because the second token (e.g. '#dot#') seems to have no content, which is not true. Maybe a solution would be for the second token to refer the first one.

### 2.1.1 Formal description of `token`

```
tok =
  element token {
    attribute id { xsd:ID }?,
    attribute value { xsd:NCName }?,
    ( attribute from { DocumentLocation },
      attribute to { DocumentLocation }
    |
      text
    )
  }
```

The information `id`, `from`, `to`, `value` are implemented as attributes of element `token`. Content of attribute `id` should be an unique identifier with XML type `ID`.

## 2.2 Word Forms as linguistic units

The segment identified by `token` elements are used to anchor word forms, that may generally be associated to a lexical entry in an lexicon and also characterized by a part of speech as well as morphological and grammatical properties.

A token may be associated to more than one word form and, conversely, a word form may cover more than one token.

For instance, in French, the morphological agglutination of *auquel* («of whose») may have several representations, depending of the properties of the segmentation:

A. The character sequence *auquel* is not decomposed and is covered by a single `token`, with two word forms covering this segment.

```
<token value="auquel" id="t0">auquel</token>
<wordForm entry="à" tag="pos@prep" tokens="t0"/>
<wordForm entry="lequel" tag="pos@pronrel" tokens="t0"/>
```

B. the segmentation identifies two agglutinated parts materialized by two tokens, each of them anchoring a word form:

```
<token value="à" id="t0">auquel</token>
<token value="lequel" id="t1"/>
<wordForm entry="à" tag="pos@prep" tokens="t0"/>
<wordForm entry="lequel" tag="pos@pronrel" tokens="t1"/>
```

**QUESTION:** Should we add the other interpretations with poly-categories

These various interpretations can be motivated by the usage or by the available tools for a given language. All of them can be described following the current recommendation.

As mentioned before, there is no linguistic mandatory properties defining the tokens, which can for instance by automatically recognized by regular languages. On the other hand, a word form, that may cover zero, one or more tokens, should represent a linguistic unit that may carry morpho-syntactic information.

The current proposal does not discuss the linguistic choices that define these linguistic units but provide enough flexibility to annotate them. The choice may be motivated by lexical or morphological properties based on context and language (depending on the *nature* and *function* of words).

### 2.2.1 Token attachment

The simplest case of relationship between tokens and word forms is when a word form covers a single token.

```
    <token id="t0" value="apple">apple</token>
2   <wordForm entry="apple" tokens="t0"/>
```

However, the current proposal allows the handling of more complex cases, as the identification of compound words covering several adjacent tokens:

```
    <token id="t0" value="prime">prime</token>
2   <token id="t1" value="minister">minister</token>
    <wordForm entry="prime_minister" tokens="t0 t1"/>
```

A sequence of non connected tokens may also be attached to a word form, for instance to handle cases where some material is inserted inside the components of a word form:

```
  <token value="afin" id="t1">afin</token>
2 <token value="justement" id="t2">justement</token>
  <token value="de" id="t3">de</token>

4
  <wordForm entry="afin_de" tokens="t1 t3"/>
6 <wordForm entry="justement" tokens="t2"/>
```

This kind of phenomena may also occur for verbs with detached particles, for instance in English or German. The English infinitive verbal form "to <verb>" may also fit in this scheme.

```
  <token value="to" id="t1">to</token>
2 <token value="eventually" id="t2">eventually</token>
  <token value="decide" id="t3">decide</token>

4
  <wordForm entry="decide" tokens="t1 t3"/>
6 <wordForm entry="eventually" tokens="t2"/>
```

In order to identify discontinuous word-form while preserving some information about the position of each component in the flow of word forms, one may use word forms covering the same sequence tokens and referring to the same entry (but possibly sub-entries).

```
  <token value="to" id="t1">to</token>
2 <token value="eventually" id="t2">eventually</token>
  <token value="decide" id="t3">decide</token>

4
  <wordForm entry="lexicon:decide:to" tokens="t1 t3"/>
6 <wordForm entry="eventually" tokens="t2"/>
  <wordForm entry="lexicon:decide:main" tokens="t1 t3"/>
```

**QUESTION:** Should be add something more explicit to signal the fact the two word forms are actually strongly related, for instance an ID/IDREF mechanism on `wordForm`.

**META: The notion of sub-entry is to be precised in the next subsection**

Another case that may arise is when one wishes to materialize a word form without any realization in the original document and, therefore, associated with an empty sequence of tokens. It may be for instance the case of some pronouns in Spanish or the hypothesis of traces.

```xml
<token value="Jean" id="t1">Jean</token>
<token value="propose" id="t2">propose</token>
<token value="de" id="t3">de</token>
<token value="partir" id="t4">partir</token>

<wordForm entry="Jean" tokens="t1"/>
<wordForm entry="propose" tokens="t2"/>
<wordForm entry="de" tokens="t3"/>
<wordForm entry="PRO"/>
<wordForm entry="partir" tokens="t4"/>
```

It should be noted that a word form is associated to some possibly empty sequence of tokens but that it is positioned relatively to other word forms and not tokens, as will be detailed in Section 2.3. It implies that there is no problem with word forms that are not associated with tokens.

Finally, several word forms may be attached to a same token:

```xml
<!-- (Donne-le moi) -->
<token value="damelo" id="t1">Damelo</token>
<wordForm entry="da" tokens="t1"/>   <!-- (Donne) -->
<wordForm entry="me" tokens="t1"/>   <!-- (le) -->
<wordForm entry="lo" tokens="t1"/>   <!-- (moi) -->
```

### 2.2.2 Relationship to lexicons

**META: Relationships to lexicons and limits: word form with entries, compound word forms**

A word form is a linguistic unit identified by its morpho-syntactic properties. Generally, this linguistic unit may be characterized by a label corresponding to an entry if some lexicon. This identification is materialized by the attribute `entry`, whose content should express a reference to the lexicon entry.

```xml
<token value="prime" id="t1">Prime</token>
<token value="minister" id="t2">minister</token>
<wordForm entry="lexicon:prime_minister" tokens="t1 t2"/>
```

The notion of lexicon entry is outside the scope of this proposal. A reference to a lexicon entry is therefore not precisely defined and, in first approximation, should correspond to some XPath pointer. In particular, one may wish reference to "sub-entries" in lexicons, for polysemous entries or for compound forms:

```xml
<token value="to" id="t1">to</token>
<token value="eventually" id="t2">eventually</token>
<token value="decide" id="t3">decide</token>

<wordForm entry="lexicon:decide:to" tokens="t1 t3"/>
<wordForm entry="eventually" tokens="t2"/>
<wordForm entry="lexicon:decide:main" tokens="t1 t3"/>
```

**META: The notation for referring lexicon entries should be precised and used in examples (XPath notation ?)**

13

However, a sequence of tokens may sometimes be identified as forming a word form because of various properties but can not associated to some lexicon entry, either becase no lexicon is available or because the word form corresponds to a named entities (proper nouns, dates, ...) or to a neologism. In that case, the content of attribute **entry** may be left empty of filled by some informative label, that is not a reference:

```
   <token id="t0">123</token>
2  <wordForm entry="NUMBER" tokens="t0"/>
```

The structure of compound forms may be expressed using nested word forms, therefore providing information about the subparts even when none is available for the whole, for instance for neologisms:

```
   <!-- birthday gift wrapping paper -->
2  <token value="Geburtstags" id="t0">Geburtstags</token>
   <token value="Geschenk" id="t1">geschenk</token>
4  <token value="Papier" id="t2">papier</token>
   <wordForm tokens="t1␣t2␣t3">
6    <wordForm entry="germanlex:geburstag" tokens="t1"/>
     <wordForm entry="germanlex:geschenk" tokens="t2"/>
8    <wordForm entry="germanlex:papier" tokens="t3"/>
   </wordForm>
```

### 2.2.3 Morpho-Syntactic content

At this point, a word form has been presented as covering a (possibly non connected) sequence of tokens, (possibly) referring to some lexicon entry. A word form should actually be completed by morpho-syntactic information precising its linguistic nature and its grammatical function in its current context.

More precisely, a *word form* carries a (complex) structure expressed by a *Feature Structure* that may attach one or more (possibly complex) values to a grammatical category or to some other linguistic property (i.e. part of speech, lemma, lexeme, morphology, ...).

**META: Exact reference of FSR**

This structure may also provides information of interest about morpho-syntax but not characterizing a word form, such as the origin of the information (a specific lexicon, a morphological tool, ...) or a confidence (returned by a stochastic tagger).

```
<token   id="t0" value="belle">belle</token>
2 <wordForm entry="lexicon:beau" tokens="t0">
  <fs>
4   <f name="cat"><sym value="adj"/></f>
    <f name="adj_type"><sym value="qual"/></f>
6   <f name="gender"><sym value="fem"/></f>
    <f name="num"><sym value="sing"/></f>
8 </fs>
  </wordForm>
```

Based on the FSR proposition, compact representation of feature structures can be used to fill the attribute **tag** of element **wordForm**.

```
<token   id="t0" value="belle">belle</token>
```

```
2  <wordForm entry="lexicon:beau" tag="cat@adj␣adj_type@qual␣gender@fem␣
        num@sing" tokens="t0"/>
```

This alternate notation will be sufficient for most cases of morpho-syntactic annotations, associated with the definition of a *tagset*.

The notion of tagset and the expressive power of Feature Structure is detailed in Section 3.

### 2.2.4  Formal description: `wordForm`

```
   wordForm =
2    element wordForm {
        attribute entry { xsd:NMTOKEN },
4       ( attribute tag { string } | fs ),
        ( attribute tokens  { xsd:IDREFS } | tok+ ),
6       wordForm*
     }
```

## 2.3  Handling ambiguities

The grouping of tokens into word forms as well as the lexical entry attached to a word form or its attached morpho-syntactic information may be ambiguous. It is therefore necessary to be able to handle these different cases of ambiguities.

A general and very generic answer is to materialize the possible readings as possible pathes through an Directed Acyclic Graph (DAG) whose edges are labeled by word form. Such DAGs forms a subpart of Finite State Automata and also cover the notion of *word lattice* used in the parsing and speech recognition communities. They can powerful enough to represent ambiguities between several decomposition into compound forms. They can also be used to denote simpler cases of lexical ambiguities.

For instance, the French textual sequence "*fer à cheval*" (horse shoe) can still be decomposed into several readings («[horse shoe]», «[iron] [on horse]», «[iron] [of] [horse]»), giving the following DAG:



Figure 2: DAG de fer à cheval

```
                    1/3 - fer - fer à cheval          2/3 - à - fer à cheval          3/3 - cheval - fer à cheval
   <token value="fer" id="t1">fer</token>
2  <token value="à" id="t2">à</token>
   <token value="cheval" id="t3">cheval</token>            1/1 - à
                                                1/1 - fer
4  <state id="S0" type="initial"/>
   <state id="S2"/>

                                                   1/2 - à - à cheval          2/2 - cheval - à cheval
```

15

```
6  <state id="S3" type="final"/>
   <transition source="S0" target ="S3">
8      <wordForm entry="lexicon:fer_à_cheval" tokens="t1␣t2␣t3"/>
   </transition >
10 <transition source="S0" target ="S1">
       <wordForm entry="lexicon:fer" tokens="t1"/>
12 </transition >
   <transition source="S1" target ="S2">
14     <wordForm entry="lexicon:à" tokens="t2"/>
   </transition >
16 <transition source="S2" target ="S3">
       <wordForm entry="lexicon:cheval" tokens="t3"/>
18 </transition >
   <transition source="S1" target ="S3">
20     <wordForm entry="lexicon:à_cheval" tokens="t2␣t3"/>
   </transition >
```

The linguistic units (entry) "fer à cheval", "fer", "à", "cheval" et "à cheval" correspond to minimal syntagmatic units that can be annotated.

Additional information could be added to edges such as probabilities.

### 2.3.1 Non ambiguous linear composition

When there is no ambiguity and one wishes to use simpler notation, the current proposal suggest a very simple alternate linear notation where element wordForm are implicitly chained following their appearance order, as illustrated by the following example:

```
   <token value="fer" id="t1">fer</token>
2  <token value="à" id="t2">à</token>
   <token value="cheval" id="t3">cheval</token>
4  <wordForm entry="fer" tokens="t1"/>
   <wordForm entry="à" tokens="t2"/>
6  <wordForm entry="cheval" tokens="t3"/>
```

This notation would be a shortcut for the following canonical version:

```
   <fsm init="s0" final="s3">
2    <token id="t0" value="fer">fer</token>
     <token id="t1" value="à">à</token>
4    <token id="t2" value="cheval">cheval</token>
     <state id="s0"/>
6    <state id="s3"/>
     <transition source="s0" target="s3">
8        <wordForm entry="lexicon:fer_à_cheval" tag="cat@noun␣..." tokens="
             t0␣t1␣t2"/>
     </transition>
10   <state id="s1"/>
     <transition source="s0" target="s1">
12       <wordForm entry="lexicon:fer" tag="cat@noun␣..." tokens="t0"/>
     </transition>
14   <transition source="s1" target="s3">
         <wordForm entry="lexicon:à_cheval" tag="cat@adv" tokens="t1␣t2"/>
16   </transition>
     <state id="s2"/>
```
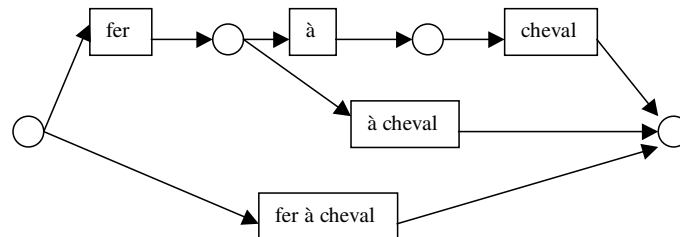
```xml
18    <transition source="s1" target="s2">
          <wordForm entry="lexicon:à" tag="cat@prep" tokens="t1"/>
20    </transition>
      <transition source="s2" target="s3">
22        <wordForm entry="lexicon:cheval" tag="cat@noun ..." tokens="t2"/>
      </transition>
24 </fsm>
```

### 2.3.2 Lexical Ambiguities

Ambiguities between different lexical entries for a same sequence of tokens can be handled by the generic element **alt**:

```xml
<token   value="porte" id="t0">porte</token>
2 <alt>
    <wordForm tokens="t0" entry="lexicon:porte"   tag="cat@noun ..."/>
4   <wordForm tokens="t0" entry="lexicon:porter" tag="cat@verb ..."/>
  </alt>
```

Such a form is actually a shortcut for the expanded notation:

```xml
<token   id="t0" value="porte">porte</token>
2 <fsm init="s0" final="s1">
  <state id="s0"/>
4 <state id="s1"/>
  <transition source="s0" target="s1">
6   <wordForm tokens="t0" entry="lexicon:porte"   tag="cat@noun ..."/>
  </transition>
8 <transition source="s0" target="s1">
    <wordForm tokens="t0" entry="lexicon:porter" tag="cat@verb ..."/>
10 </transition>
  </fsm>
```

### 2.3.3 Word form Content Ambiguities

The proposal on Feature Structure Representation provides several ways to represent ambiguities, for instance at the level of feature values. These mechanisms may be imported to avoid the use to DAG for ambiguities within the morpho-syntactic content of a word-form.

For instance, French verb form "mange" (to eat) is ambiguous between 1st and 3rd persons, which can be represented by the **vAlt** element proposed in FSR.

```xml
  <token   id="t0" value="mange">mange</token>
2 <wordForm tokens="t0" entry="lexicon:manger">
      <fs>
4       <f name="cat"><sym value="verb"/></f>
        <f name="aux"><sym value="avoir"/></f>
6       <f name="mode"><sym value="indicatif"/></f>
        <f name="tense"><sym value="present"/></f>
8       <f name="pers">
            <vAlt>
10              <sym value="1"/>
                <sym value="3"/>
```

```
12                    </vAlt>
                </f>
14              <f name="num"><sym value="sing"/></f>
            </fs>
16      </wordForm>
```

A compact notation can still be used by registering most frequent cases of ambiguities in Feature Structure Libraries (see section about Tagsets).

```
  <token   id="t0" value="mange">mange</token>
2 <wordForm tokens="t0" entry="lexicon:manger" tag="cat@verb␣aux@avoir␣
      mode@ind␣tense@pres␣pers@1|3␣num@sing"/>
```

### 2.3.4   Mixing simple and complex cases

Ambiguities are generally localized and it is tempting to also localize the use of the graph notation only where it is needed. The current proposal provides a way to switch between linear and graph notations.

```
  <token   id="t0" value="afin">afin</token>
2 <token id="t1" value="de">de</token>
  <fsm init="s0" final="s2">
4   <transition source="s0" target="s2">
        <wordForm entry="lexicon:afin_de" tag="cat@prep" tokens="t0␣t1"/>
6   </transition>
    <transition source="s0" target="s1">
8       <wordForm entry="lexicon:afin" tag="cat@prep" tokens="t0"/>
    </transition>
10    <transition source="s1" target="s2">
        <wordForm entry="lexicon:de" tag="cat@prep" tokens="t1"/>
12    </transition>
  </fsm>
14 <token id="t2" value="grandir">grandir</token>
  <wordForm entry="lexicon:grandir" tag="cat@verb␣..." tokens="t2"/>
16 <token id="t3" value=",">,</token>
  <wordForm entry="lexicon:," tag="cat@ponct" tokens="t3"/>
18 <token id="t4" value="il">il</token>
  <wordForm entry="lexicon:il" tag="cat@pronoun␣..." tokens="t4"/>
20 <token id="t5" value="mange">mange</token>
  <wordForm entry="lexicon:manger" tag="cat@verb␣..." tokens="t5"/>
22 <token id="t6" value="des(1:2)">des</token>
  <wordForm entry="des(1:2)" tag="cat@adv" tokens="t6"/>
24 <token id="t7" value="des(2:2)"/>
  <wordForm entry="lexicon:le" tag="cat@det␣..." tokens="t7"/>
26 <token id="t8" value="pommes">pommes</token>
  <token id="t9" value="de">de</token>
28 <token id="t10" value="terre">terre</token>
  <fsm init="s8" final="s11">
30    <transition source="s8" target="s11">
        <wordForm entry="lexicon:pomme_de_terre" tag="cat@noun␣..." tokens=
            "t8␣t9␣t10"/>
32    </transition>
    <transition source="s8" target="s9">
```

```
34        <wordForm entry="lexicon:pomme" tag="cat@noun ⌣..." tokens="t8"/>
      </transition>
36    <transition source="s9" target="s10">
          <wordForm entry="lexicon:de" tag="cat@prep" tokens="t9"/>
38    </transition>
      <transition source="s10" target="s11">
40        <wordForm entry="lexicon:terre" tag="cat@noun ⌣..." tokens="t10"/>
      </transition>
42 </fsm>
```

A shortcut notation for the rather long expanded version:

```
   <fsm init="s0" final="s11">
2    <token id="t0" value="afin">afin</token>
     <token id="t1" value="de">de</token>
4    <state id="s0"/>
     <state id="s2"/>
6    <transition source="s0" target="s2">
         <wordForm entry="lexicon:afin_de" tag="cat@prep" tokens="t0⌣t1"/>
8    </transition>
     <state id="s1"/>
10   <transition source="s0" target="s1">
         <wordForm entry="lexicon:afin" tag="cat@prep" tokens="t0"/>
12   </transition>
     <transition source="s1" target="s2">
14       <wordForm entry="lexicon:de" tag="cat@prep" tokens="t1"/>
     </transition>
16   <token id="t2" value="grandir">grandir</token>
     <state id="s3"/>
18   <transition source="s2" target="s3">
         <wordForm entry="lexicon:grandir" tag="cat@verb ⌣..." tokens="t2"/>
20   </transition>
     <token id="t3" value=",">,</token>
22   <state id="s4"/>
     <transition source="s3" target="s4">
24     <wordForm entry="lexicon:," tag="cat@ponct" tokens="t3"/>
     </transition>
26   <token id="t4" value="il">il</token>
     <state id="s5"/>
28   <transition source="s4" target="s5">
         <wordForm entry="lexicon:il" tag="cat@pronoun ⌣..." tokens="t4"/>
30   </transition>
     <token id="t5" value="mange">mange</token>
32   <state id="s6"/>
     <transition source="s5" target="s6">
34     <wordForm entry="lexicon:manger" tag="cat@verb ⌣..." tokens="t5"/>
     </transition>
36   <token id="t6" value="des(1:2)">des</token>
     <state id="s7"/>
38   <transition source="s6" target="s7">
         <wordForm entry="des(1:2)" tag="cat@adv" tokens="t6"/>
40   </transition>
     <token id="t7" value="des(2:2)"/>
42   <state id="s8"/>
```

```
   <transition source="s7" target="s8">
44      <wordForm entry="lexicon:le" tag="cat@det ⌴..." tokens="t7"/>
   </transition>
46 <token id="t8" value="pommes">pommes</token>
   <token id="t9" value="de">de</token>
48 <token id="t10" value="terre">terre</token>
   <state id="s11"/>
50 <transition source="s8" target="s11">
      <wordForm entry="lexicon:pomme_de_terre" tag="cat@noun ⌴..." tokens="
         t8⌴t9⌴t10"/>
52 </transition>
   <state id="s9"/>
54 <transition source="s8" target="s9">
      <wordForm entry="lexicon:pomme" tag="cat@noun ⌴..." tokens="t8"/>
56 </transition>
   <state id="s10"/>
58 <transition source="s9" target="s10">
      <wordForm entry="lexicon:de" tag="cat@prep" tokens="t9"/>
60 </transition>
   <transition source="s10" target="s11">
62      <wordForm entry="lexicon:terre" tag="cat@noun ⌴..." tokens="t10"/>
   </transition>
64 </fsm>
```

## 2.4 Formal description: `fsm`, `state`, `transition`

```
   fsm =
2    element fsm {
     attribute init {xsd:IDREF}?,
4    attribute final {xsd:IDREF}?,
     (tok | state | transition )+
6    }
   state =
8    element state {
             attribute id { xsd:ID }
10   }
   transition =
12   element transition {
         attribute source { xsd:IDREF },
14       attribute target { xsd:IDREF },
         (wordForm | wordFormAlt)
16   }
   wordFormAlt =
18   element alt { wordForm+ }
```

Also a question about alternatives on word forms: should we use a generic element `alt` or define a specialized version `wordFormAlt`.

20

# 3 Morpho-Syntactic annotation content

The previous section explain how to complete a document with morpho-syntactic annotations. However, it does not precise the content of these annotations. What set of features and feature values to use to express this content (within element `wordForm`) and with which meaning ?

Such a set is usually referred as a *tagset* specifying the content of possible annotations. However, the diversity of approaches and languages makes almost impossible the proposition of an unique tagset. More modestly or pragmatically, the current proposal is to provide mechanisms to define tagsets by relying on a Data Category Registry (DCR) and Feature Structures (FSR).

An annotated document will therefore be completed by either adding or referring to a *tagset*.

## 3.1 Representing Feature Structures

This section illustrates, trough a few examples, how the FSR proposal may be used to represent morpho-syntactic content. A morpho-syntactic content attached to a word form may be represented by a set of pairs (attribute, value) called *feature specifications*. Following FSR, such a simple set is represented by:

```
   <wordForm entry="prime_minister" tokens="t1_t2">
2  <fs>
         <f name="pos">
4           <sym value="noun"/>
         </f>
6        <f name="gender">
             <sym value="masculine"/>
8        </f>
         <f name="number">
10          <sym value="singular"/>
         </f>
12 </fs>
   </wordForm>
```

### 3.1.1 Disjunctive values

FSR provides way to denotes alternatives for values.

```
   <token   value="mange" id="t0">mange</token>
2  <wordForm entry="lexicon:manger" tokens="t0">
     <fs>
4      <f name="cat"><sym value="verb"/></f>
       <f name="aux"><sym value="avoir"/></f>
6      <f name="mode"><sym value="indicatif"/></f>
       <f name="tense"><sym value="present"/></f>
8      <f name="pers">
         <vAlt>
10          <sym value="1"/>
            <sym value="3"/>
12       </vAlt>
```

```
            </f>
14        <f name="num"><sym value="sing"/></f>
        </fs>
16 </wordForm>
```

### 3.1.2  Multiple Tags

One may wish to assign a sequence of tags to a word form, for instance for some contracted form. A way to do that is to use the multiple values `vMult` provided by FSR.

```
   <token   id="t0" value="auxquels">auxquels</token>
2 <wordForm entry="lexicon:auquel" tokens="t0">
     <vMult org="list">
4        <fs>
             <f name="cat"><sym value="prep"/></f>
6        </fs>
         <fs>
8            <f name="cat"><sym value="pronoun"/></f>
             <f name="kind"><sym value="rel"/></f>
10           <f name="num"><sym value="pl"/></f>
             <f name="gender"><sym value="masc"/></f>
12       </fs>
     </vMult>
14 </wordForm>
```

Note: In general, an alternate solution for contracted word is to set up several word forms attached to a same token, each of them with a simple tag, even if each word form is some subpart of a lexical entry.

```
   <token   id="t0" value="auxquels">auxquels</token>
2 <wordForm entry="lexicon:auquel/à" tokens="t0">
       <fs>
4            <f name="cat"><sym value="prep"/></f>
       </fs>
6 </wordForm>
  <wordForm   entry="lexicon:auquel/quel" tokens="t0">
8      <fs>
             <f name="cat"><sym value="pronoun"/></f>
10           <f name="kind"><sym value="rel"/></f>
             <f name="num"><sym value="pl"/></f>
12           <f name="gender"><sym value="masc"/></f>
       </fs>
14 </wordForm>
```

## 3.2  Compact notations for morpho-syntactic annotations

FSR proposal provides ways for compact representations of feature structure, essentially by defining *libraries* of feature structures, feature specifications, and feature values. These mechanisms may be used to fill the content of attribute `tag` for element `wordForm` by compact tags, following a standard practice in the NLP community.

The generic way provided by FSR to use libraries is illustrated by the following example, with the attribute `feats` of element `fs`:

```
    <!-- A feature value library -->
2  <vLib>
        <sym id="noun" value="noun"/>
4       <sym id="sing" value="singular"/>
        <sym id="masc" value="masculine"/>
6  </vLib>
    <!-- A feature specification library -->
8  <fvLib>
        <f id="pos@n" name="pos" fVal="noun"/>
10      <f id="num@sing" name="num" fVal="sing"/>
        <f id="gen@fem" name="gen" fVal="fem"/>
12 </fvLib>
    <!-- Annotated document -->
14 <wordForm entry="prime_minister" tokens="t1">
        <fs feats ="pos@n_num@sing_gen@fem"/>
16 </wordForm>
```

The current proposal allows the use of attribute `tag` to assign a compact tag directly to a word form:

```
<wordForm entry="prime_minister" tokens="t1" tag="pos@n_num@sing_
    gen@masc"/>
```

Even disjunctive values may be simplified, following the same mechanism:

```
    <!-- A feature value library -->
2  <vLib>
        <vAlt id="1|3">
4           <sym value="1"/>
            <sym value="3"/>
6       </vAlt>
        ...
8  </vLib>
    <!-- A feature specification library -->
10 <fvLib>
        <f id="pers@1|3" name="pers" fVal="1|3"/>
12      ...
    </fvLib>
14 <!-- Annotated document -->
    <token id="t0" value="porte">porte</token>
16 <wordForm entry="lexicon:porter" tag="cat@verb_pers@1|3_num@sing_..."
        tokens="t0"/>
```

## 3.3 Designing tagsets with Morpho-Syntactic Data Categories

**META: This section should be completed.**

The attributes, values, and possibly feature types used to specify morpho-syntactic content are not just labels but carry linguistic meanings, in other word a semantic. To avoid misinterpretations, this semantic should be clearly defined.

However, the current proposal does not try to define the semantic of an unique complete set of such attributes, values, and types. It would be almost impossible to be complete

(given the diversity of languages) and almost impossible to assign to each element a semantic accepted by the whole community.

Instead, it is proposed that an annotated document should be completed by the indication of one or more *tagsets*.

The first objective of a tagset is to list the terminology used to annotate a document, precising its semantic w.r.t, an official registered terminology, namely a *Data Category Registery* for Morpho-Syntax. The process may be seen as selecting a set of data categories (*Data Category Selection – DCS*).

A private category may be introduced by selecting some registered one:

```
  <dcs private="genre" registered="dcs:morphosyntax:gender:french" rel="
      eq"/>
2 <dcs private="fem" registered="
      dcs:morphosyntax:gender:french:femininin"/>
```

The correspondance with a registered category may not be perfect. The `rel` may be used to specify which relationship exists between the private and registered categories. For instance, one may introduce a private category advneg as being subsumed by a more general registered category adverb.

```
  <dcs private="advneg" registered="dcs:morphosyntax:pos:adverb" rel="
      subs"/>
2 <dcs private="strange" rel="none"/>
```

It is also possible (but not advised) to introduce a private category bearing no relationship with any registered category.

```
  <dcs private="title"/>
```

When the correspondance is not perfect, a few words of description should be added to precise the meaning of the private category.

```
  <dcs private="title">
2     <description> A part of speech used to denote honorific titles like
      Pr. or S.A.S.
4     </description>
  </dcs>
```

**QUESTION:** The proposed scheme only cover very elementary mappings. Should we propose more sophisticated mappings as equivalence between feature structures,

The second objective of an tagset is to specify the set of valid feature structures based on the introduced terminology. It will be achieved by relying on a ISO standard to come about Feature Structure Declaration.

A third objective of an tagset is to name the most common morphosyntactic structures through the use of feature structure libraries, in order to be able to use compact tags for annotations.

### 3.3.1 Formal description: `tagset`

**META: TO BE DEFINED**

```
tagset =
2    element tagset {
```

```
            (  attribute  ref  {  xsd:anyURL  }
4            |  dcs ∗ , fsd , fsLibs ∗
            )
6        }
   dcs =
8       element  dcs  {
            attribute  private  {  xsd:NCName  } ,
10          (  attribute  registered  {  xsd:anyURL  } ,
              attribute  rel  {  "eq"  |  "subs"  |  "gen"  }  ) ?
12          element  description  {  text  } ?
       }
14 ## fsLibs :  to  be  imported
   ## fsd :  to  be  imported  ( and  defined )
```

The `dcs` corresponds to a Data Category Selection part whose exact content is still to be defined.

The `fsd` corresponds to a Feature Structure Declaration part whose normalization is yet to be done.

# 4   Global information

The current proposal does not address the issue of adding metadata, globally and on annotations, for instance to mention the date (creation, revision) of an annotation and its author(s).

A global element `msa` (for **Morpho-Syntactic Annotations**) is introduced as root element to contain morpho-syntactic annotations. This element `msa` is used to specify global properties relative to the annotated documents, in particular the addressing scheme used to delimit segments. One or more tagsets should also be mentioned either explicitly or by reference with an URL.

**META: TO BE COMPLETED**

```
<msa  addressing ="mpeg7">
2   <token  id="t0"  from=""  to=""  value=""">
  ␣␣<wordForm␣tokens="t0"␣entry="" >␣…␣</wordForm>
4 ␣␣…
  </msa>
```

## 4.1   Formal description

```
start =
2       element  maf  {
          addressing ,
4         tagset ,
          ( fsm  |  tok  |  wordForm  |  wordFormAlt  )+
6       }
   ## To  be  precised
8 addressing =
      attribute  addressing  {  xsd:NMTOKEN  }
10 DocumentLocation  =  xsd:NMTOKEN
```

25

META: The notion of *addressing scheme reference* is yet to be precised. A possible list of such schema would include:

- TEI ptrs,

- XPath,

- byte offsets,

- Unicode char offsets,

- MPEG7 multimedia addressing

# A   Informative Relax NG compact schema

```
  # $Id$
2
  # Preliminary Relax NG schema for MAF: Morpho Syntactic Annotations
4 # Eric de la Clergerie <Eric.De_La_Clergerie@inria.fr>

6 default namespace = ""

8 start =
    element maf {
10       addressing ,
         tagset ,
12       (fsm | tok | wordForm | wordFormAlt )+
    }
14 ##To be precised
  addressing =
16    attribute addressing { xsd:NMTOKEN }
  tok =
18    element token {
         attribute id { xsd:ID }?,
20       attribute value { xsd:NCName }?,
         ( attribute from { DocumentLocation },
22         attribute to { DocumentLocation }
          |
24         xsd:NMTOKEN
         )
26    }
  wordForm =
28    element wordForm {
         attribute entry { xsd:NMTOKEN },
30       ( attribute tag { string } | fs ),
         ( attribute tokens  { xsd:IDREFS } | tok+ ),
32       wordForm*
    }
34 fsm =
    element fsm {
36       attribute init {xsd:IDREF}?,
         attribute final {xsd:IDREF}?,
38       (tok | state | transition )+
    }
40 state =
    element state {
42       attribute id { xsd:ID }
    }
44 transition =
    element transition {
46       attribute source { xsd:IDREF },
         attribute target { xsd:IDREF },
48       (wordForm | wordFormAlt)
    }
50 wordFormAlt =
    element alt { wordForm+ }
```

27

```
52 tagset =
       element tagset {
54         ( attribute ref { xsd:anyURL }
            | dcs∗,fsd ,fsLibs∗
56        )
       }
58 dcs =
       element dcs {
60        attribute private { xsd:NCName },
          ( attribute registered { xsd:anyURL },
62           attribute rel { "eq" | "subs" | "gen" } )?
          element description { text }?
64     }
   ## To be precised
66 DocumentLocation = xsd:NMTOKEN
```

## B   Informative DTD

```
  <?xml version="1.0" encoding="UTF−8"?>
2

4    <!−−

6    DTD for Morpho−Syntaxic Annotation Framework

8   −−>

10 <!ENTITY % fs −SYSTEM  'INCLUDE'>
   <!ENTITY % fs −PUBLIC  'IGNORE'>
12
   <![% fs −SYSTEM [
14        <!ENTITY % dtd−fs SYSTEM "fs.dtd">
   ]]>
16 <![% fs −PUBLIC [
          <!ENTITY % dtd−fs PUBLIC "−//DTD fs //DTD//EN" "http://www.
             tc37sc4.org/dtd/fs.dtd">
18 ]]>
   %dtd−fs;
20

22   <!ELEMENT msa          ((token | wordForm | wordFormAlt | fsm)+) >

24   <!ELEMENT token        (#PCDATA) >
     <!ATTLIST token        value CDATA #IMPLIED
26                          from IDREF #IMPLIED
                            to IDREF #IMPLIED
28                          id ID #IMPLIED>

30   <!ELEMENT wordForm      ( fs | wordForm | token)∗>
     <!ATTLIST wordForm      tokens IDREFS #IMPLIED
32                          entry CDATA #IMPLIED
```

```
                                  tag CDATA #IMPLIED>
34

36    <!ELEMENT wordFormAlt ( wordForm+)>

38    <!ELEMENT fsm         (( state | transition | token)+)>
      <!ATTLIST fsm         initial IDREF #REQUIRED
40                          finals  IDREFS #REQUIRED>

42    <!ELEMENT state       EMPTY>
      <!ATTLIST state       id ID #REQUIRED>
44
      <!ELEMENT transition  (wordFormAlt | wordForm)>
46    <!ATTLIST transition  from IDREF #REQUIRED
                            to IDREF #REQUIRED>
```