

FoLiA: Format for Linguistic Annotation

version 0.8 - revision 2.2

Documentation

ILK Technical Report – ILK 12-03

Maarten van Gompel
Center for Language Studies
Radboud University Nijmegen

P.O. Box 9103, NL-6500 HD Nijmegen, The Netherlands

URL: <http://ilk.uvt.nl/clam/>¹

February 23, 2012

¹This document is available from <http://ilk.uvt.nl/downloads/pub/papers/ilk.1203.pdf>.
All rights reserved, Tilburg University and Radboud University Nijmegen.

Contents

1	Introduction	4
2	Document Format	8
2.1	Global Structure	8
2.2	Identifiers	9
2.3	Basic Structural Elements	9
2.4	Paradigm & Terminology	13
2.5	Annotation Declarations	15
2.6	Structure Annotation	16
2.6.1	Divisions	16
2.6.2	Gaps	17
2.6.3	Whitespace and Linebreaks	19
2.6.4	Events	19
2.6.5	Lists	20
2.6.6	Figures	21
2.6.7	Tables	21
2.7	Token Annotation	21

2.7.1	Part of Speech Annotation	22
2.7.2	Lemma Annotation	22
2.7.3	Phonetic Annotation	23
2.7.4	Lexical Semantic Sense Annotation	23
2.7.5	Domain Tags	24
2.7.6	Corrections	24
2.8	Alternative Token Annotations	32
2.9	Span Annotation	33
2.9.1	Entities	33
2.9.2	Syntax	34
2.9.3	Dependency Relations	35
2.9.4	Chunking	38
2.9.5	Events as span annotation	39
2.9.6	Semantic roles	40
2.9.7	Alterative Span Annotations	40
2.10	Advanced Paradigm	41
2.10.1	Human readable Descriptions	41
2.10.2	Text content and multiple classes	42
2.10.3	Features, Subsets and multiple classes	45
2.10.4	Part-of-Speech tags with features	47
2.11	Subtoken Annotation	48
2.11.1	Morphological Analysis	48

2.11.2	Named Entities within a token	49
2.12	Alignments	49
2.13	Metadata	51
3	Set Definition Format	53
3.1	Introduction	53
3.2	Types and classes	53
3.3	Concept link	54
3.4	Subsets	54
3.5	Constraints	55
A	Common Queries	59
B	Validation	61
C	Implementations	62
D	Conversions	64
E	Comparison	66
E.0.1	TCF v0.4 – Text Corpus Format	66
E.0.2	KAF – Kyoto Annotation Format	67
E.0.3	LAF – Linguistic Annotation Framework, and GraF	67
E.0.4	TEI P5 – Text Encoding Initiative	67

Chapter 1

Introduction

FoLiA is a Format for Linguistic Annotation, derived from the D-Coi format[2] developed as part of the D-Coi project by project partner at Polderland Language and Speech Technologies B.V. The D-Coi format was designed for use by the D-Coi corpus, as well as by its successor; the SoNaR corpus [6]. Though being rooted in the D-Coi format, the FoLiA format goes a lot further and introduces a rich generalised framework for linguistic annotation. FoLiA development started at the ILK research group, Tilburg University, and is continued at the Radboud University Nijmegen. It is being adopted in multiple projects in the Dutch and Flemish Natural Language Processing community.

FoLiA is an XML-based[4] annotation format, suitable for the representation of written and linguistically annotated language resources. FoLiA's intended use is as a format for storing and/or exchanging language resources, including corpora. Our aim is to introduce a single rich format that can accomodate a wide variety of linguistic annotation types through a single generalised paradigm. We do not commit to any label set, language or linguistic theory. This is always left to the developer of the language resource, and provides maximum flexibility.

XML is an inherently hierarchic format. FoLiA does justice to this by maximally utilising a hierarchic, inline, setup. We inherit from the D-Coi format, which posits to be loosely based on a minimal subset of TEI[5]. Because of the introduction of a new and much broader paradigm, FoLiA is *not* backwards-compatible with D-Coi, i.e. validators for D-Coi will not accept FoLiA XML. It is however easy to convert FoLiA to less complex or verbose formats such as the D-Coi format, or plain-text. Converters will be provided. This may entail some loss of information if the simpler format has no provisions for particular types of

information specified in the FoLiA format.

The most important characteristics of FoLiA are:

- **Generalised** paradigm - We use a generalised paradigm, with as few ad-hoc provisions for annotation types as possible.
- **Expressivity** - The format is highly expressive, annotations can be expressed in great detail and with flexibility to the user's needs. Moreover, FoLiA has generalised support for representing annotation alternatives, and annotation metadata such as information on annotator, time of annotation, and annotation confidence.
- **Extensible** - Due to the generalised paradigm and the fact that the format does not commit to any label set, FoLiA is fairly easily extensible.
- **Formalised** - The format is formalised, validatable on both a shallow and a deep level (the latter including tagset validation), and easily machine parsable, for which tools are provided.
- **Practical** - FoLiA has been developed in a bottom-up fashion right alongside applications, libraries, and other toolkits and convertors. Whilst the format is rich, we try to maintain it as simple and straightforward as possible, minimising the learning curve and making it easy to adopt FoLiA in practical applications.

The FoLiA format makes mixed-use of inline and stand-off annotation. Inline annotation is used for annotations pertaining to single tokens, whilst stand-off annotation in a separate annotation layers is adopted for annotation types that span over multiple token. This provides FoLiA with the necessary flexibility and extensibility to deal with various kinds of annotations. Inspiration for this was in part obtained from the Kyoto Annotation Format [1].

In publication of research that makes use of FoLiA, a citation should be given of: *"Maarten van Gompel (2012). FoLiA: Format for Linguistic Annotation. Documentation. ILK Technical Report 12-03.*

Available from <http://ilk.uvt.nl/downloadcds/pub/papers/ilk.1203.pdf>".

FoLiA is open-source and all technical resources are licensed under the GNU Public License v3.

Further notable features of the FoLiA format include:

- XML-based, validation against RelaxNG schema.
- Full Unicode support; UTF-8 encoded.
- Document structure consists of divisions, paragraphs, sentences and words/tokens, and more specific elements.
- Can be used for both tokenised as well as untokenised text, though for meaningful linguistic annotation, tokenisation is mandatory.
- Provenance support for all linguistic annotations: annotator, type (automatic or manual), time.
- Support for alternative annotations, optionally with associated confidence values.
- Support for features using subsets, allowing for more detailed user-defined annotation.
- Not committed to any label set, these are user-defined.
- Agnostic with regard to metadata. External metadata schemes such as CMDI are recommended.

There is support for the following linguistic annotations:

- Part-of-Speech tags (with features)
- Lemmatisation
- Spelling corrections on both a tokenised as well as an untokenised level
- Lexical semantic sense annotation (to be used in DutchSemCor)
- Named Entities / Multi-word units
- Syntactic Parses
- Dependency Relations
- Chunking
- Morphological Analysis
- Subjectivity Annotation (to be used in VU-DNC)

In later stages, the following may be added:

- Semantic Role Labelling
- Co-reference
- Topic Segmentation
- Authorship Attribution

FoLiA support is incorporated directly into the following ILK software:

- ucto - A tokeniser which can directly output FoLiA XML
- Frog - A PoS-tagger/lemmatiser/parser suite (the successor of Tadpole), will eventually support reading and writing FoLiA.
- CLAM - Computational Linguistics Application Mediator, will eventually have viewers for the FoLiA format.
- PyNLPI - Python Natural Language Processing Library, comes with a library for parsing FoLiA
- libfolia - C++ library for parsing FoLiA

FoLiA is used in the following projects (list may not be complete):

- SoNaR
- DutchSemCor (NWO)
- TTNWW (CLARIN)
- Valkuil.net
- CLAM
- Frog
- Ucto

To clearly understand this documentation, note that when we speak of “elements” or “attributes”, we refer to XML notation, i.e. XML elements and XML attributes.

Chapter 2

Document Format

2.1 Global Structure

In FoLiA, each document/text is represented by one XML file. The basic structure of such a FoLiA document is as follows and should always be UTF-8 encoded. An elaborate XSLT stylesheet will be provided in order to be able to instantly view FoLiA documents in any modern web browser.

```
<?xml version="1.0" encoding="utf-8"?>
<?xml-stylesheet type="text/xsl" href="http://ilk.uvt.nl/FoLiA/FoLiA.xsl"?>
<FoLiA xmlns="http://ilk.uvt.nl/FoLiA"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  version="0.5"
  xml:id="example">
  <metadata>
    <annotations>
      ...
    </annotations>
    <!-- (Here CMDI or IMDI metadata can be inserted) -->
  </metadata>
  <text xml:id="example.text">
    ...
  </text>
</FoLiA>
```

2.2 Identifiers

Many elements in the FoLiA format specify an identifier by which the element is uniquely identifiable. This makes referring to any part of a FoLiA document easy and follows the lead of the D-Coi format. Identifiers should be unique in the entire document, and can be anything that qualifies as a valid ID according to the XML standard. A well proven convention is of a cumulative nature, in which you append the elements name, a period, and a sequence number, to the identifier of a parent element higher in the hierarchy. Identifiers are always encoded in the `xml:id` attribute. The FoLiA document as a whole also carries an ID.

Identifiers are very important and used throughout the FoLiA format. They enable external resources and database to easily point to a specific part of the document or its annotation. FoLiA has been set up in such a way that *identifiers should never ever change*. Once an identifier is assigned, it should never change, re-numbering is strictly prohibited unless you intentionally want to create a new resource and break compatibility with the old one.

2.3 Basic Structural Elements

Basic structural elements occur within the `text` element. These are the most basic ones:

- `p` - Paragraph
- `s` - Sentence
- `w` - Word (token)

These are typically nested, the word elements cover the actual tokens. This is the most basic level of annotation; tokenisation. Let's take a look at an example, we have the following text:

This is a paragraph containing only one sentence.

This is the second paragraph. This one has two sentences.

In FoLiA XML, this will appear as follows after tokenisation. Some parts have been omitted for the sake of brevity:

```
<p xml:id="example.p.1">
  <s xml:id="example.p.1.s.1">
    <w xml:id="example.p.1.s.1.w.1"><t>This</t></w>
    <w xml:id="example.p.1.s.1.w.2"><t>is</t></w>
    ...
    <w xml:id="example.p.1.s.1.w.8" space="no"><t>sentence</t></w>
    <w xml:id="example.p.1.s.1.w.9"><t>.</t></w>
  </s>
</p>
<p xml:id="example.p.2">
  <s xml:id="example.p.2.s.1">
    <w xml:id="example.p.2.s.1.w.1"><t>This</t></w>
    <w xml:id="example.p.2.s.1.w.2"><t>is</t></w>
    ..
    <w xml:id="example.p.2.s.1.w.5" space="no"><t>paragraph</t></w>
    <w xml:id="example.p.2.s.1.w.6"><t>.</t></w>
  </s>
  <s xml:id="example.p.2.s.2">
    <w xml:id="example.p.2.s.2.w.1"><t>This</t></w>
    <w xml:id="example.p.2.s.2.w.2"><t>one</t></w>
    ..
    <w xml:id="example.p.2.s.2.w.5" space="no"><t>sentences</t></w>
    <w xml:id="example.p.2.s.2.w.6"><t>.</t></w>
  </s>
</p>
```

The deepest content element should always contain a text element (t) which holds the actual textual content. The necessity of having a text element shall become apparent as you progress through this documentation; there can be many different token annotations under a word element (w).

FoLiA is not just a format for holding tokenised text, although tokenisation is a prerequisite for almost all kinds of annotation. However, FoLiA can also hold untokenised text, on a paragraph and/or sentence level:

```
<p xml:id="example.p.1">
  <s xml:id="example.p.1.s.1">
    <t>This is a paragraph containing only one sentence.</t>
  </s>
</p>
<p xml:id="example.p.2">
  <s xml:id="example.p.2.s.1">
    <t>This is the second paragraph.</t>
  </s>
  <s xml:id="example.p.2.s.2">
```

```

        <t>This one has two sentences.</t>
    </s>
</p>

```

Higher level elements *may* also contain a text element even when the deeper element does too. It is very important to realise that the sentence/paragraph-level text element *always* contains the text *prior* to tokenisation! Note also that the word element has an attribute `space`, which defaults to `yes`, and indicates whether the word was followed by a space in the *untokenised* original. This allows for partial reconstructibility of the sentence in its untokenised form. See section 2.10.2 for a more elaborate overview of this subject.

The following example shows the maximum amount of redundancy, with text elements at every level.

```

<p xml:id="example.p.1">
    <t>This is a paragraph containing only one sentence.</t>
    <s xml:id="example.p.1.s.1">
        <t>This is a paragraph containing only one sentence.</t>
        <w xml:id="example.p.1.s.1.w.1"><t>This</t></w>
        <w xml:id="example.p.1.s.1.w.2"><t>is</t></w>
        ...
        <w xml:id="example.p.1.s.1.w.8" space="no"><t>sentence</t></w>
        <w xml:id="example.p.1.s.1.w.9"><t>.</t></w>
    </s>
</p>
<p xml:id="example.p.2">
    <t>This is the second paragraph. This one has two sentences.</t>
    <s xml:id="example.p.2.s.1">
        <t>This is the second paragraph.</t>
        <w xml:id="example.p.2.s.1.w.1"><t>This</t></w>
        <w xml:id="example.p.2.s.1.w.2"><t>is</t></w>
        ..
        <w xml:id="example.p.2.s.1.w.5" space="no"><t>paragraph</t></w>
        <w xml:id="example.p.2.s.1.w.6"><t>.</t></w>
    </s>
    <s xml:id="example.p.2.s.2">
        <t>This one has two sentences.</t>
        <w xml:id="example.p.2.s.2.w.1"><t>This</t></w>
        <w xml:id="example.p.2.s.2.w.2"><t>one</t></w>
        ..
        <w xml:id="example.p.2.s.2.w.5" space="no"><t>sentences</t></w>
        <w xml:id="example.p.2.s.2.w.6"><t>.</t></w>
    </s>
</p>

```

It this kind of redundancy is used (it is not mandatory),you may optionally point back to the text content of its parent by specifying the offset attribute:

```
<p xml:id="example.p.1">
  <t>This is a paragraph containing only one sentence.</t>
  <s xml:id="example.p.1.s.1">
    <t offset="0">This is a paragraph containing only one sentence.</t>
    <w xml:id="example.p.1.s.1.w.1"><t offset="0">This</t></w>
    <w xml:id="example.p.1.s.1.w.2"><t offset="5">is</t></w>
    ...
    <w xml:id="example.p.1.s.1.w.8" space="no"><t offset="40">sentence</t></w>
    <w xml:id="example.p.1.s.1.w.9"><t offset="48">.</t></w>
  </s>
</p>
```

Matters can become more complicated as multiple text-content element of different classes may be associated with an element, this will be discussed later on in section 2.10.2.

Paragraph elements may be omitted if a document is described that does not distinguish paragraphs but only sentences. Sentences however may never be omitted; FoLiA documents can never consist of tokens only!

The content element head is reserved for headers and captions, it behaves similarly to the paragraph element and may hold sentences.

FoLiA also explicitly supports quotes, as demonstrated in the next example, which annotates the following sentence:

He said: ‘‘I do not know . I think you are right. ", and left.

A quote may consist of one or more sentences, but may also consist of mere tokens. The token identifiers in all cases simply follow the sequential numbering of the root sentence, not the embedded sentence.

```
<s xml:id="example.p.1.s.1">
  <w xml:id="example.p.1.s.1.w.1" class="WORD"><t>He</t></w>
  <w xml:id="example.p.1.s.1.w.2" class="WORD"><t>said</t></w>
  <w xml:id="example.p.1.s.1.w.3" class="PUNCTUATION" space="no"><t>:</t></w>
  <w xml:id="example.p.1.s.1.w.4" class="PUNCTUATION" space="no"><t>' '</t></w>
  <quote xml:id="example.p.1.s.1.quote.1">
    <s xml:id="example.p.1.s.1.quote.1.s.1">
      <w xml:id="example.p.1.s.1.w.5" class="WORD"><t>I</t></w>
      <w xml:id="example.p.1.s.1.w.6" class="WORD"><t>do</t></w>
      <w xml:id="example.p.1.s.1.w.7" class="WORD"><t>not</t></w>
```

```

    <w xml:id="example.p.1.s.1.w.8" class="WORD"><t>know</t></w>
    <w xml:id="example.p.1.s.1.w.9" class="PUNCTUATION" space="no"><t>.</t></w>
  </s>
  <s xml:id="example.p.1.s.1.quote.1.s.2">
    <w xml:id="example.p.1.s.1.w.10" class="WORD"><t>I</t></w>
    <w xml:id="example.p.1.s.1.w.11" class="WORD"><t>think</t></w>
    <w xml:id="example.p.1.s.1.w.12" class="WORD"><t>you</t></w>
    <w xml:id="example.p.1.s.1.w.13" class="WORD"><t>are</t></w>
    <w xml:id="example.p.1.s.1.w.14" class="WORD"><t>right</t></w>
  </s>
</quote>
<w xml:id="example.p.1.s.1.w.15" class="PUNCTUATION" space="no"><t>'</t></w>
<w xml:id="example.p.1.s.1.w.16" class="PUNCTUATION"><t>,</t></w>
<w xml:id="example.p.1.s.1.w.17" class="WORD"><t>and</t></w>
<w xml:id="example.p.1.s.1.w.18" class="WORD"><t>left</t></w>
<w xml:id="example.p.1.s.1.w.19" class="PUNCTUATION" space="no"><t>.</t></w>
</s>

```

2.4 Paradigm & Terminology

The FoLiA format has a very uniform setup and its XML notation for annotation follows a generalised paradigm. We distinguish three different categories of annotation, of which the latter two apply to actual linguistic annotation:

- **Structural annotation** - Annotations marking global structure, such as chapters, sections, subsections, figures, list items, paragraphs, etc...
- **Token annotation** - Annotations pertaining to one specific token. These will be elements of the token element (`w`) in inline notation. Linguistic annotations in this category are for example: part-of-speech annotation, lemma annotation, sense annotation, morphological analysis, spelling correction. Some token elements may be used on higher levels (e.g. sentence/paragraph) as well and may then be referred to as **Extended Token Annotation**
- **Span annotation** - Annotations spanning over multiple tokens. Each type of annotation will be in a separate **annotation layer** with offset notation. These layers are typically embedded on the sentence level, or possibly also on higher levels (paragraph/division/text) for certain annotation types. Examples in this category are: syntactic parses, chunking, semantic roles and named entities.

- **Subtoken annotation** - This is a type of annotation that has aspects of both token annotation as well as span annotation. The former because is included inline within a token element (`w`) and describes the token. The latter because it defines a small span *within* the token. Examples in this category are: morphological analysis, named subentities

Almost all annotations are associated with what we shall call a **set**. The set determines the vocabulary, the tags or types, of the annotation. An element of such a set is referred to as a **class** from the FoLiA perspective. For example, we may have a document with Part-of-Speech annotation according to the CGN set (a tagset for Dutch part-of-speech tags). The CGN set defines main tag classes such as *WW*, *BW*, *ADJ*, *VZ*. FoLiA itself thus never commits to any tagset but leaves you to explicitly define this. You can also use multiple tagsets in the same document if so desired, even for the same type of annotation.

Any annotation element may have a `set` attribute, the value of which points to the URL hosting the file that defines the set, and a `class` attribute, which selects a class from the set.

The following example shows a simple Part-of-Speech annotation without features, but with all common attributes according to the FoLiA paradigm:

```
<pos set="http://ilk.uvt.nl/folia/sets/CGN" class="WW"
  annotator="Maarten_van_Gompel" annotortype="manual"
  confidence="0.76" datetime="1982-12-15T19:01" />
```

The example demonstrates that any annotation element can take an `annotator` attribute and an `annotortype`. The latter is either “manual” for human annotators, or “auto” for automated systems. The value for `annotator` is open and should be set to the name or ID of the system or human annotator that made the annotation. The `confidence` attribute can be set to a floating point value between zero and one, and expresses the confidence the annotator places in his annotation. Last, the `datetime` attribute specifies the date and time when this annotation was recorded, the format is YYYY-MM-DDThh:mm:ss (note the T in the middle to separate date from time), as per the XSD Datetime data type. None of these options are mandatory, only `class` may be mandatory for some types of annotation, such as `pos`.

More advanced aspects of the paradigm will be introduced later in section 2.10.

2.5 Annotation Declarations

The annotation declaration is a mandatory part of the metadata that declares all the types of annotation and all sets that are present in the document. Annotations are declared in the `annotations` block, as shown in the following example. We here define four annotation levels.

```
<annotations>
  <token-annotation set="http://ilk.uvt.nl/olia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
  <pos-annotation set="http://ilk.uvt.nl/olia/sets/CGN"
    annotator="Frog" annotortype="auto" />
  <lemma-annotation set="http://ilk.uvt.nl/olia/sets/lemmas-nl"
    annotator="Frog" annotortype="auto" />
  <sense-annotation set="http://ilk.uvt.nl/olia/sets/Cornetto"
    annotator="SupWSD1" annotortype="auto" />
</annotations>
```

The `set` attribute is mandatory¹ and refers to a URL of a FoLiA Set Definition file (see chapter ??). The Set Definition specifies exactly what classes are allowed in the set. It for examples specifies exactly what Part-of-Speech tags exists. This information is necessary to completely validate the document at its deepest level. If the sets point to URLs that do not exist, warnings will be issued. Validation can still proceed but with the notable exception of deep validation of these sets.

If multiple sets are used for the same annotation type, they each need a separate declaration:

```
<pos-annotation set="http://ilk.uvt.nl/olia/sets/CGN"
  annotator="Frog" annotortype="auto" />
<pos-annotation set="http://ilk.uvt.nl/olia/sets/brown" />
```

If only one set is declared, then in the document itself you are allowed to skip the `set` attribute on these specific annotation elements. The declared set will automatically be the default.

The `annotator` and `annotortype` attributes act as defaults, for the specific annotation type and set. Unlike `set`, you do *not* need, and it is in fact prohibited, to declare every possible annotator here!

Annotator defaults can always be overridden at the specific annotation elements. But declaring them allows for the annotation element to be less verbosely ex-

¹Technically, it can be omitted, but then the set defaults to “undefined”. This is allowed for flexibility and less explicit usage of FoLiA in limited settings, but not recommended!

pressed. Explicitly referring to a set and annotator for each annotation element can be cumbersome and pointless in a document with a single set and a single annotator for that particular type of annotation. Declarations and defaults provide a nice way around this problem.

2.6 Structure Annotation

2.6.1 Divisions

Within the text element, the structure element `div` can be used to create divisions and subdivisions. Each division *may* be of a particular *class* pertaining to a *set* defining all possible classes. Divisions and other structural units are often numbered, think for example of chapters and sections. The number, as it was in the source document, can be encoded in the `n` attribute of the structure annotation element.

Look at the following example, showing a full FoLiA document with structured divisions:

```
<?xml version="1.0" encoding="utf-8"?>
<?xml-stylesheet type="text/xsl" href="http://ilk.uvt.nl/FoLiA/FoLiA.xsl"?>
<FoLiA xmlns="http://ilk.uvt.nl/FoLiA"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  version="0.5"
  xml:id="example">
  <metadata>
    <annotations>
      <div-annotation set="http://ilk.uvt.nl/fofia/sets/divisions" />
    </annotations>
    <!-- (Here CMDI or IMDI metadata can be inserted) -->
  </metadata>
  <text xml:id="example.text">
    <div class="chapter" n="1">
      <head><t>Introduction</t></head>
      <div class="section" n="1">
        <div class="subsection" n="1.1">
          <t>In the beginning....</t>
        </div>
      </div>
      ...
    </div>
  </text>
</FoLiA>
```

If divisions with sets are present, they need to be declared, as can be seen in the metadata. If you have divisions without sets or classes, then a declaration is not necessary

Divisions stem from D-Coi and are modified in FoLiA. These divisions are not mandatory, but may be used to mark extra structure. D-Coi supported the elements `div0`, `div1`, `div2`, etc., but FoLiA only knows a single `div` element, which can be nested at will and associated with classes. Note that paragraphs, sentences and words have their own explicit tags, as seen earlier, divisions should never be used for marking these, only larger structures can be divisions!

The `head` element may be used to for the header of any division. It may hold `s` and `w` elements (not `p`).

2.6.2 Gaps

Sometimes there are parts of a document you want to skip and not annotate, but include as is. For this purpose the `gap` element should be used. Gaps may have a particular class indicating the kind of gap it is. Common omissions are for example front-matter and back-matter.

The D-Coi format pre-defined the following “reasons” [2]:

- `frontmatter`
- `backmatter`
- `illegible`
- `other-language`
- `cancelled`
- `inaudible`
- `sampling`

Due to the flexible nature of FoLiA, we don’t predefine any classes whatsoever and leave this up to whatever set is declared. The above gives a good indication of what gaps can be used for though.

The `gap` element may optionally take two elements:

1. desc - holding a substitute that may be shown to the user, describing what has been omitted.
2. content - The actual raw content of the omission, as it was without further annotations. This is an XML CDATA type element, excluding it from any kind of parsing.

```

<text xml:id="example.text">
  <gap class="frontmatter" annotator="Maarten_van_Gompel">
    <desc>This is the cover of the book</desc>
    <content>
<![CDATA[
    SHOW WHITE AND THE SEVEN DWARFS

    by the Brothers Grimm

    first edition

    Copyright(c) blah blah
]]>
  </content>
</gap>
<div class="chapter" n="1">
  <head><t>Introduction</t></head>
  <div class="section" n="1">
    <div class="subsection" n="1.1">
      <t>In the beginning....</t>
    </div>
  </div>
  ...
</div>
</text>

```

Gaps have to be declared:

```

<annotations>
  <gap-annotation set="http://ilk.uvt.nl/foia/sets/dcoi-gaps" />
</annotations>

```

2.6.3 Whitespace and Linebreaks

Sometimes you may want to explicitly specify vertical whitespace or line breaks. This can be done using respectively `whitespace` and `br`. Both are simple structural elements that need not be declared. Note that using `p` to denote paragraphs is always strongly preferred over using `whitespace` to mark their boundaries!

```
<text xml:id="example.text">
  <s xml:id="example.s.1">
    <w xml:id="example.s.1.w.1">
      <br />
    <w xml:id="example.s.1.w.2">
      <w xml:id="example.s.1.w.3">
    </s>
  <whitespace />
  <s xml:id="example.s.2">
  </s>
</text>
```

The difference between `br` and `whitespace` is that the former specifies that only a linebreak was present, not forcing any vertical whitespace, whilst the latter actually generates an empty space, which would be comparable to two successive `br` statements. Both elements can be used inside divisions, paragraphs, and sentences.

2.6.4 Events

Event structure, though uncommon to regular written text, can be useful in certain documents. Divisions, paragraphs, sentences, or even words can be encapsulated in an event element to indicate they somehow form an event entity of a particular class. This kind of structure annotation is especially useful in dealing with written media such as chat logs, tweets, and internet fora, in which chat turns, forum posts, and tweets can be demarcated as particular events.

Below an example of a simple chat log, word tokens omitted for brevity:

```
<event class="message" begindatetime="2011-12-15T19:01"
  enddatetime="2011-12-15T19:05" actor="Jane_Doe">
  <s>
    <t>Hello John.</t>
  </s>
  <s>
    <t>How are you doing?</t>
```

```

    </s>
</event>
<event class="message" begintatetime="2011-12-15T19:06"
  actor="John_Doe">
  <s>
    <t>I am fine Jane, thanks.</t>
  </s>
</event>

```

The (optional) features `begintatetime` and `enddatetime` can be used express the exact moment at which an event started or ended. Note that this differs from the generic `datetime` attribute, which would describe the time at which the annotation was recorded, rather than when the event took place! Also, `begintatetime` and `enddatetime` are so-called *features* (see section

The declaration:

```

<annotations>
  <event-annotation set="http://ilk.uvt.nl/olia/sets/events" />
</annotations>

```

For more fine-grained control over timed events, for example within sentences. It is recommended to use the `timedevent` span annotation element instead! This works in a very similar fashion but uses an offset annotation layer. See section 2.9.5.

2.6.5 Lists

FoLiA, like D-Coi, allows lists to be explicitly marked as shown in the following example:

```

<head><t>My grocery list</t></head>
<list xml:id="example.list.1">
  <item xml:id="example.list.1.item.1" n="A"><t>Apples</t></item>
  <item xml:id="example.list.1.item.2" n="B"><t>Pears</t></item>
</list>

```

The `item` element may hold sentences (`s`) and words (`w`). The D-Coi format had a `label` element, this is deprecated in favour of the `n` attribute in the `item` itself.

Lists, like paragraphs, sentences and headers are content elements that need not be declared and are not associated with a set or class.

2.6.6 Figures

Even figures can be encoded in the FoLiA format, although the actual figure itself can only be included as a mere reference to an external image file, but including such a reference (`src` attribute) is optional.

```
<figure xml:id="example.figure.1" n="1" src="/path/or/url/to/image/file">
  <desc>A textual description of the figure (Like ALT in HTML)</desc>
  <caption><t>The caption for the figure</t></caption>
</figure>
```

Figures are not declared. The `caption` element may hold sentences (`s`) and words (`w`).

2.6.7 Tables

Development Notes
There is no provision yet for encoding tables. This may be added later.

2.7 Token Annotation

Token annotations are annotations that are placed within the word (`w`) element. They all can take any of the attributes described in section 2.4, this has to be kept in mind when reading this section. Moreover, all token annotations depend on the document being tokenised, i.e. there being `w` elements. If the tokens are assigned classes part of a certain set, then a declaration such as the following is mandatory:

```
<annotations>
  <token-annotation set="http://ilk.uvt.nl/folia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
</annotations>
```

Being part of a set, this implies that tokens themselves *may* be assigned a class, as is for example done by the tokeniser *ucto*:

```
<s xml:id="example.p.1.s.1">
  <w xml:id="example.p.1.s.1.w.1" class="WORD"><t>I</t></w>
  <w xml:id="example.p.1.s.1.w.2" class="WORD"><t>see</t></w>
```

```

    <w xml:id="example.p.1.s.1.w.3" class="NUMBER">2</w>
    <w xml:id="example.p.1.s.1.w.4" class="WORD" space="no">children</w>
    <w xml:id="example.p.1.s.1.w.5" class="PUNCTUATION">.</w>
</s>

```

2.7.1 Part of Speech Annotation

The following example illustrates a simple Part-of-Speech annotation for the Dutch word “boot”:

```

<w xml:id="example.p.1.s.1.w.2">
  <t>boot</t>
  <pos class="N" />
</w>

```

However, for some tagsets simple part-of-speech annotation is not enough; there may for example be features associated with the part of speech tag. We will into this later, in section 2.10.

Whenever Part-of-Speech annotations are used, they should be declared in the annotations block as follows, the set you use may differ and all attributes are optional. In the declaration example here it is as if the annotations were made by the software *Frog*. Do note the requirement of a token-annotation as well.

```

<annotations>
  <token-annotation set="http://ilk.uvt.nl/foia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
  <pos-annotation set="http://ilk.uvt.nl/foia/sets/CGN"
    annotator="Frog" annotortype="auto" />
</annotations>

```

As mentioned earlier, the declaration only sets defaults. They can be overridden in the pos element itself (or any other token annotation element for that matter).

2.7.2 Lemma Annotation

In the FoLiA paradigm, lemmas are perceived as classes within the (possibly open) set of all possible lemmas. Their annotation is thus as follows:

```

<w xml:id="example.p.1.s.1.w.2">
  <t>boot</t>

```

```

    <lemma class="boot" />
</w>

```

And the example declaration:

```

<annotations>
  <token-annotation set="http://ilk.uvt.nl/foia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
  <lemma-annotation set="http://ilk.uvt.nl/foia/sets/mblem-nl"
    annotator="Frog" annotortype="auto" />
</annotations>

```

2.7.3 Phonetic Annotation

Phonetic annotations can be included as follows. Similarly to lemmas, they may often refer to a set with possibly open classes.

```

<w xml:id="example.p.1.s.1.w.2">
  <t>boot</t>
  <phon-annotation set="http://ilk.uvt.nl/foia/sets/ipa"
    class="bu:t" />
</w>

```

This is an extended token annotation element that can also be used directly on a sentence or paragraph level.

And the example declaration:

```

<annotations>
  <token-annotation set="http://ilk.uvt.nl/foia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
  <phon-annotation set="http://ilk.uvt.nl/foia/sets/ipa" />
</annotations>

```

2.7.4 Lexical Semantic Sense Annotation

In semantic sense annotation, the classes in most sets will be a kind of lexical unit ID. In systems that make a distinction between lexical units and synonym sets (synsets), the synset attribute is available for notation of the latter. In systems with only synsets and no other primary form of lexical unit, the class can simply be set to the synset.

A human readable description for the *sense* element, “beeldhouwwerk”, can be placed inside a *desc* element, but this is optional.

```
<w xml:id="example.p.1.s.1.w.2">
  <t>beeld</t>
  <sense class="r_n-6220" synset="d_n-32683"><desc>beeldhouwwerk</desc></sense>
</w>
```

The example declaration is as follows:

```
<annotations>
  <token-annotation set="http://ilk.uvt.nl/foia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
  <sense-annotation set="http://ilk.uvt.nl/foia/sets/cornetto" />
</annotations>
```

2.7.5 Domain Tags

This is an extended token annotation element, which means it can also be used directly in any of the content elements, such as sentence (s) and paragraph (p). It can even be used in the *text* element itself. This annotation defines the domain of the token of content element. Example:

```
<w xml:id="example.p.1.s.1.w.2">
  <t>boot</t>
  <domain class="nautical" />
</w>
```

The declaration:

```
<annotations>
  <token-annotation set="http://ilk.uvt.nl/foia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
  <domain-annotation set="http://ilk.uvt.nl/foia/sets/domains-nl" />
</annotations>
```

2.7.6 Corrections

Corrections, including but not limited to spelling corrections, can be annotated using the *correction* element. It can be applied as an extended token annotation element as in the following example, which shows a spelling correction of the misspelled word “tree” to its corrected form “tree”.

```

<w xml:id="example.p.1.s.1.w.1">
  <correction xml:id="TEST-000000001.p.1.s.1.w.1.c.1" class="spelling">
    <new>
      <t>tree</t>
    </new>
    <original auth="no">
      <t>treee</t>
    </original>
  </correction>
</w>

```

The class indicates the kind of correction, according to the set used. The new elements holds the actual content of the correction. The original element holds the content prior to correction. Note that all corrections must carry a unique identifier. In this example, what we are correcting is the actual textual content, the text element (t). To facilitate the job of parsers and queriers, the original element has to be marked as being non-authoritative, using auth="no". This states that this element and anything below it is not authoritative, meaning that any text or annotations within do not affect the text or annotations of the structure element (the word in this case) of which it is a part.

Corrections can be nested and we want to retain a full back-log. The following example illustrates the word "treee" that has been first mis-corrected to "three" and subsequently corrected again to "tree":

```

<w xml:id="example.p.1.s.1.w.1">
  <correction xml:id="TEST-000000001.p.1.s.1.w.1.c.2" class="spelling"
    annotator="Jane_Doe" annotortype="manual" confidence="1.0">
    <new>
      <t>tree</t>
    </new>
    <original auth="no">
      <correction xml:id="TEST-000000001.p.1.s.1.w.1.c.1" class="spelling"
        annotator="John_Doe" annotortype="manual" confidence="0.6">
        <new>
          <t>three</t>
        </new>
        <original auth="no">
          <t>treee</t>
        </original>
      </correction>
    </original>
  </correction>
</w>

```

In the examples above what we corrected was the actual textual content (t). It is however also possible to correct other annotations: The next example corrects a part-of-speech tag; in such cases, there is no t element in the correction, but simply another token annotation element, or group thereof.

```
<w xml:id="example.p.1.s.1.w.1">
  <t>tree</t>
  <correction xml:id="TEST-000000001.p.1.s.1.w.1.c.1">
    <new>
      <pos class="n" />
    </new>
    <original auth="no">
      <pos class="v" />
    </original>
  </correction>
</w>
```

Error detection and correction with suggestions

The correction of an error implies the detection of an error. In some cases, detection comes without correction, for instance when the generation of correction suggestions is postponed to a later processing stage. The `errordetection` element is a very simple element that serves this purpose. It signals the existence of errors, or absence thereof:

```
<w xml:id="example.p.1.s.1.w.1">
  <t>treee</t>
  <errordetection class="spelling" annotator="errorlistX" error="yes" />
</w>
```

The `error` attribute is set to “yes” (which is the default value), and thus marks this as an error of class “spelling”. We can also imagine it specifically marking something as *not* being an error (in which case `class` is always redundant), for example due to the occurrence of the word according to a lexicon:

```
<w xml:id="example.p.1.s.1.w.1">
  <t>tree</t>
  <errordetection annotator="lexiconX" error="no" />
</w>
```

This kind of error detection is very simple and does not provide actual correction nor suggestions for correction. In some cases, it is desirable to record suggestions for correction, but without making the actual correction.

The correction tag can also be used in such situations in which you want to list suggestions for correction, but not yet commit to any single one. You may for example want to postpone this actual selection to another module or human annotator. Recall that the actual correction is always included in the “new” tag, non-committing suggestions are included in the “suggestion” tag. All suggestions may take an ID and may specify an annotator, if no annotator is specified it will be inherited from the correction element itself. Suggestions never take sets or classes by themselves, the class and set pertain to the correction as a whole, and apply to all suggestions within. This implies that you will need multiple correction elements if you want to make suggestions of very distinct types. The following example shows two suggestions for correction:

```
<w xml:id="example.p.1.s.1.w.1">
  <t>treee</t>
  <correction xml:id="example.p.1.s.1.w.1.c.1"
    class="spelling" annotator="errorlistX">
    <suggestion confidence="0.8" auth="no">
      <t>tree</t>
    </suggestion>
    <suggestion confidence="0.2" auth="no">
      <t>three</t>
    </suggestion>
  </correction>
</w>
```

In the situation above we have a possible correction with two suggestions, none of which has been selected yet. The actual text remains unmodified so there are no new or original tags. Note that anything in the scope of a suggestion is by definition non-authoritative and suggestions have to be marked as such using `auth="no"` to facilitate the job of parsers.

When an actual correction is made, the correction element changes. It may still retain the list of suggestions. In the following example, a human annotator named John Doe took one of the suggestions and made the actual correction:

```
<w xml:id="example.p.1.s.1.w.1">
  <correction xml:id="example.p.1.s.1.w.1.c.1"
    class="spelling" annotator="John_Doe"
    annotortype="human">
    <new>
      <t>tree</t>
    </new>
    <suggestion annotator="errorlistX" auth="no"
      annotortype="auto" confidence="0.8">
      <t>tree</t>
    </suggestion>
  </correction>
</w>
```

```

        <suggestion annotator="errorlistX" auth="no"
            annotortype="auto" confidence="0.2">
            <t>three</t>
        </suggestion>
        <original auth="no">
            <t>treee</t>
        </original>
    </correction>
</w>

```

Something similar may happen when a correction is made *on the basis of* one or more kinds of error detection, the correction element directly embeds the errordetection element:

```

<w xml:id="example.p.1.s.1.w.1">
    <correction class="spelling" annotator="John_Doe">
        <new>
            <t>tree</t>
        </new>
        <original auth="no">
            <t>treee</t>
        </original>
        <errordetection class="spelling" annotator="errorlist" annotortype="auto" />
    </correction>
</w>

```

In the above example, “treee” was detected by an automated error list as being an error, and was corrected to “tree” by human annotator John Doe.

Like everything, corrections and error detection have to be declared, and have to be declared separately. Nothing stops you from pointing them both to the same set however. Suggestions fall under the scope of corrections and need not be declared separately.

```

<annotations>
    <token-annotation set="http://ilk.uvt.nl/olia/sets/ucto-tokconfig-nl"
        annotator="ucto" annotortype="auto" />
    <errordetection-annotation set="http://ilk.uvt.nl/olia/sets/corrections" />
    <correction-annotation set="http://ilk.uvt.nl/olia/sets/corrections" />
</annotations>

```

Merges, Splits and Swaps

Sometimes, one wants to merge multiple tokens into one single new token, or the other way around; split one token into multiple new ones. The FoLiA format does

not allow you to simply create new tokens and reassign identifiers. Identifiers are by definition permanent and should never change, as this would break backward compatibility. So such a change is therefore by definition a correction, and one uses the correction tag to merge and split tokens.

We will first demonstrate a merge of two tokens (“on line”) into one (“online”), the original tokens are always retained as w-original elements. First a peek at the XML prior to merging:

```
<s xml:id="example.p.1.s.1">
  <w xml:id="example.p.1.s.1.w.1">
    <t>on</t>
  </w>
  <w xml:id="example.p.1.s.1.w.2">
    <t>line</t>
  </w>
</s>
```

And after merging:

```
<s xml:id="example.p.1.s.1">
  <correction xml:id="example.p.1.s.1.c.1" class="merge">
    <new>
      <w xml:id="example.p.1.s.1.w.1-2">
        <t>online</t>
      </w>
    </new>
    <original auth="no">
      <w xml:id="example.p.1.s.1.w.1">
        <t>on</t>
      </w>
      <w xml:id="example.p.1.s.1.w.2">
        <t>line</t>
      </w>
    </original>
  </correction>
</s>
```

Note that the correction element, being a kind of extended token annotation, is here a member of the sentence (s), rather than the word token (w) as in all previous examples. The new identifier denotes the span of the merge, but this is mere convention.

Now we will look at a split, the reverse of the above situation. Prior to splitting, assume we have:

```
<s xml:id="example.p.1.s.1">
```

```

<w xml:id="example.p.1.s.1.w.1">
  <t>online</t>
</w>
</s>

```

After splitting:

```

<s xml:id="example.p.1.s.1">
  <correction xml:id="example.p.1.s.1.c.1" class="split">
    <new>
      <w xml:id="example.p.1.s.1.w.1_1">
        <t>on</t>
      </w>
      <w xml:id="example.p.1.s.1.w.1_2">
        <t>line</t>
      </w>
    </new>
    <original auth="no">
      <w xml:id="example.p.1.s.1.w.1">
        <t>online</t>
      </w>
    </original>
  </correction>
</s>

```

The same principle as used for merges and splits can also be used for performing “swap” corrections:

```

<s xml:id="example.p.1.s.1">
  <correction xml:id="example.p.1.s.1.c.1" class="split">
    <new>
      <w xml:id="example.p.1.s.1.w.2">
        <t>on</t>
      </w>
      <w xml:id="example.p.1.s.1.w.1">
        <t>line</t>
      </w>
    </new>
    <original auth="no">
      <w xml:id="example.p.1.s.1.w.1">
        <t>line</t>
      </w>
      <w xml:id="example.p.1.s.1.w.2">
        <t>on</t>
      </w>
    </original>
  </correction>

```

</s>

Note that in such a swap situation, the identifiers of the word tokens will appear out of sequence after correction, due to the principle that identifiers never change once set.

Omissions and Insertions

Omissions, words that are removed in correction, and insertions, words inserted during correction, are dealt with in a way similar to merges, splits and swaps. For omissions, the new element is simply empty. In the following example the word “the” was duplicated and removed in correction:

```
<s xml:id="example.p.1.s.1">
  <w xml:id="example.p.1.s.1.w.1">
    <t>the</t>
  </w>
  <correction xml:id="example.p.1.s.1.c.1" class="duplicate">
    <new>
      </new>
      <original auth="no">
        <w xml:id="example.p.1.s.1.w.2">
          <t>the</t>
        </w>
      </original>
    </correction>
    <w xml:id="example.p.1.s.1.w.3">
      <t>man</t>
    </w>
  </s>
```

For insertions, the original element is empty.

```
<s xml:id="example.p.1.s.1">
  <w xml:id="example.p.1.s.1.w.1">
    <t>the</t>
  </w>
  <correction xml:id="example.p.1.s.1.c.1" class="duplicate">
    <new>
      <w xml:id="example.p.1.s.1.w.1_1">
        <t>old</t>
      </w>
    </new>
    <original auth="no">
      </original>
    </correction>
  </s>
```



```

</correction>
<w xml:id="example.p.1.s.1.w.2">
  <t>man</t>
</w>
</s>

```

2.8 Alternative Token Annotations

The FoLiA format does not just allow for a single authoritative annotation per token, in addition it allows for the recording of *alternative* annotations. Alternative token annotations are grouped within one or more `alt` elements. If multiple annotations are grouped together under the same `alt` element, then they are deemed dependent and form a single set of alternatives.

Each alternative has a unique identifier, formed in the already familiar fashion. In the following example we see the Dutch word “bank” in the sense of a sofa, alternatively we see two alternative annotations with a different sense and domain. Any annotation element within an *alt* block by definition needs to be marked as non-authoritative by setting `auth="no"`. This facilitates the job of parsers and quierers.

```

<w xml:id="example.p.1.s.1.w.1">
  <t>bank</t>
  <domain class="furniture" />
  <sense class="r_n-5918" synset="d_n-21410"
    annotator="John_Doe" annotortype="manual"
    confidence="1.0">zitmeubel</sense>
  <alt xml:id="example.p.1.s.1.w.1.alt.1">
    <domain auth="no" class="finance" />
    <sense auth="no" class="r_n-5919" synset="d_n-27025"
      annotator="Jane_Doe" annotortype="manual"
      confidence="0.6">geldverlenende instelling</sense>
  </alt>
  <alt xml:id="example.p.1.s.1.w.1.alt.2">
    <domain auth="no" class="geology" />
    <sense auth="no" class="r_n-5920" synset="d_n-38257"
      annotator="Jim_Doe" annotortype="manual"
      confidence="0.1">zandbank</sense>
  </alt>
</w>

```

2.9 Span Annotation

Not all annotations can be realised as token annotations. Some typically span multiple tokens. For these we introduce a kind of offset notation in separate *annotation layers*. These annotation layers are embedded at the sentence level, *after* the word tokens. Within these layers, references are made to these word tokens. Each annotation layer is specific to a kind of span annotation.

The layer elements themselves may also take the `set`, `annotator`, `annotatortype`, or `confidence` attributes. Which introduces the defaults for all the span annotations under it. They in turn may of course always chose to override this.

2.9.1 Entities

Named entities or other multi-word units can be encoded in the entities layer. Below is an example of a full sentence in which one name is tagged. Each entity should have a unique identifier.

```
<s xml:id="example.p.1.s.1">
  <t>The Dalai Lama greeted him.</t>
  <w xml:id="example.p.1.s.1.w.1"><t>The</t></w>
  <w xml:id="example.p.1.s.1.w.2"><t>Dalai</t></w>
  <w xml:id="example.p.1.s.1.w.3"><t>Lama</t></w>
  <w xml:id="example.p.1.s.1.w.4"><t>greeted</t></w>
  <w xml:id="example.p.1.s.1.w.5"><t>him</t></w>
  <w xml:id="example.p.1.s.1.w.6"><t>.</t></w>
  <entities>
    <entity xml:id="example.p.1.s.1.entity.1" class="person">
      <wref id="example.p.1.s.1.w.2" t="Dalai" />
      <wref id="example.p.1.s.1.w.3" t="Lama" />
    </entity>
  </entities>
</s>
```

Note that elements that are not part of any span annotation need never be included in the layer. The `wref` element takes an *optional* `t` attribute which contains a copy of the text of the word pointed towards. This is to facilitate human readability and prevent the need for resolving words for simple applications in which only the textual content is of interest.

2.9.2 Syntax

A very typical form of span annotation is syntax annotation. This is done within the syntax layer and introduces a nested hierarchy of syntactic unit (su) elements. Each syntactic unit should have a unique identifier.

```
<s xml:id="example.p.1.s.1">
  <t>The Dalai Lama greeted him.</t>
  <w xml:id="example.p.1.s.1.w.1"><t>The</t></w>
  <w xml:id="example.p.1.s.1.w.2"><t>Dalai</t></w>
  <w xml:id="example.p.1.s.1.w.3"><t>Lama</t></w>
  <w xml:id="example.p.1.s.1.w.4"><t>greeted</t></w>
  <w xml:id="example.p.1.s.1.w.5"><t>him</t></w>
  <w xml:id="example.p.1.s.1.w.6"><t>.</t></w>
  <syntax>
    <su xml:id="example.p.1.s.1.su.1" class="s">
      <su xml:id="example.p.1.s.1.su.1_1" class="np">
        <su xml:id="example.p.1.s.1.su.1_1_1" class="det">
          <wref id="example.p.1.s.1.w.1" t="The" />
        </su>
        <su xml:id="example.p.1.s.1.su.1_1_2" class="pn">
          <wref id="example.p.1.s.1.w.2" t="Dalai" />
          <wref id="example.p.1.s.1.w.3" t="Lama" />
        </su>
      </su>
      <su xml:id="example.p.1.s.1.su.1_2" class="vp">
        <su xml:id="example.p.1.s.1.su.1_2_1" class="v">
          <wref id="example.p.1.s.1.w.4" t="greeted" />
        </su>
        <su xml:id="example.p.1.s.1.su.1_2_2" class="pron">
          <wref id="example.p.1.s.1.w.5" t="him" />
        </su>
      </su>
    </su>
  </syntax>
</s>
```

Just to prevent any misunderstanding, the classes depend on the set used, so you can use whatever system of syntactic annotation you desire. Moreover, any of the su elements can have the common attributes `annotator`, `annotortype` and `confidence`.

The above example illustrated a fairly simple syntactic parse. Dependency parses are possible too. Dependencies are listed separate from the syntax in an extra annotation layer, as shall be explained in the next section.

The declaration is as follows:

```
<annotations>
  <token-annotation set="http://ilk.uvt.nl/olia/sets/ucto-tokconfig-nl"
    annotator="ucto" annotortype="auto" />
  <syntax-annotation set="http://ilk.uvt.nl/olia/sets/syntax-nl" />
</annotations>
```

2.9.3 Dependency Relations

Dependency relations are relations between tokens or spans of tokens, in most cases equal to syntactic units. A dependency relation is often of a particular class and consists of a single head component and a single dependent component. In the sample “He sees”, there is syntactic dependency between the two words: “sees” is the head, and “He” is the dependant, and the relation class is something like “subject”, as the dependant is the subject of the head word. Each dependency relation is explicitly noted.

The element `dependencies` introduces this annotation layer. Within it, dependency elements describe all dependency pairs.

In the below example, we show a Dutch sentence parsed with the Alpino Parser [3]. We show not only the dependency layer, but also the syntax layer to which it is related. The dependency element always contains one head element (`hd`) and one dependant element (`dep`), both can refer to a syntactic unit (or anything else for that matter) by means of the `aref` element within their scope. Additionally, the words they cover are reiterated in the usual fashion, using `wref`. For a better understanding, the figure below illustrates the syntactic parse with the dependency relations.

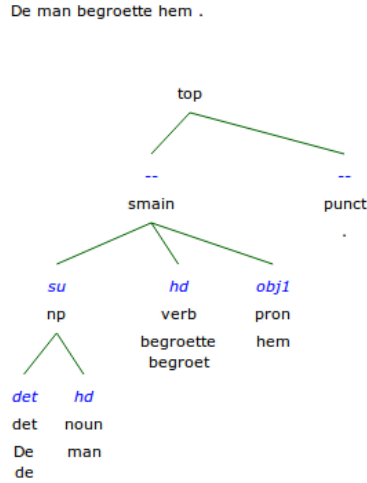


Figure 2.1: Alpino dependency parse for the Dutch sentence “De man begroette hem.”

```

<s xml:id="example.p.1.s.1">
  <t>De man begroette hem.</t>
  <w xml:id="example.p.1.s.1.w.1"><t>De</t></w>
  <w xml:id="example.p.1.s.1.w.2"><t>man</t></w>
  <w xml:id="example.p.1.s.1.w.3"><t>begroette</t></w>
  <w xml:id="example.p.1.s.1.w.4"><t>hem</t></w>
  <w xml:id="example.p.1.s.1.w.5"><t>.</t></w>
  <syntax>
    <su xml:id="example.p.1.s.1.su.1" class="top">
      <su xml:id="example.p.1.s.1.su.1_1" class="smain">
        <su xml:id="example.p.1.s.1.su.1_1_1" class="np">
          <su xml:id="example.p.1.s.1.su.1_1_1_1" class="top">
            <wref id="example.p.1.s.1.w.1" t="De" />
          </su>
          <su xml:id="example.p.1.s.1.su.1_1_1_2" class="top">
            <wref id="example.p.1.s.1.w.2" t="man" />
          </su>
        </su>
        <su xml:id="example.p.1.s.1.su.1_1_2" class="verb">
          <wref id="example.p.1.s.1.w.3" t="begroette" />
        </su>
        <su xml:id="example.p.1.s.1.su.1_1_3" class="pron">
          <wref id="example.p.1.s.1.w.4" t="hem" />
        </su>
      </su>
    </su>
  </syntax>
</s>

```

```

        <su xml:id="example.p.1.s.1.su.1_2" class="punct">
            <wref id="example.p.1.s.1.w.5" t="." />
        </su>
    </su>
</syntax>
<dependencies>
    <dependency xml:id="example.p.1.s.1.dependency.1" class="su">
        <hd>
            <wref id="example.p.1.s.1.w.3" t="begroette">
                <aref id="example.p.1.s.1.su.1_1_2" type="su">
                    </hd>
                <dep>
                    <wref id="example.p.1.s.1.w.2" t="man" />
                    <aref id="example.p.1.s.1.su.1_1_1" type="su">
                        </dep>
                    </dependency>
                <dependency xml:id="example.p.1.s.1.dependency.3" class="obj1">
                    <hd>
                        <wref id="example.p.1.s.1.w.3" t="begroette">
                            <aref id="example.p.1.s.1.su.1_1_2" type="su">
                                </hd>
                            <dep>
                                <wref id="example.p.1.s.1.w.4" t="hem" />
                                <aref id="example.p.1.s.1.su.1_1_3" type="su">
                                    </dep>
                                </dependency>
                            <dependency xml:id="example.p.1.s.1.dependency.2" class="det">
                                <hd>
                                    <wref id="example.p.1.s.1.w.2" t="man" />
                                    <aref id="example.p.1.s.1.su.1_1_1_2" type="su">
                                        </hd>
                                    <dep>
                                        <wref id="example.p.1.s.1.w.1" t="De" />
                                        <aref id="example.p.1.s.1.su.1_1_1_1" type="su">
                                            </dep>
                                        </dependency>
                                    </dependencies>
                                </s>

```

Note that in the first dependency relation, the dependent is just “man” rather than “de man” . That is, we point only to the head of dependents, the full scope follows automatically when building the dependency tree.

The declaration:

```

<annotations>
    <token-annotation set="http://ilk.uvt.nl/foolia/sets/ucto-tokconfig-nl"
        annotator="ucto" annotortype="auto" />

```

```

    <syntax-annotation set="http://ilk.uvt.nl/foia/sets/alpino-syntax" />
    <dependency-annotation set="http://ilk.uvt.nl/foia/sets/alpino-dep" />
</annotations>

```

2.9.4 Chunking

Unlike a full syntactic parse, chunking is not nested. The layer for this type of linguistic annotation is predictably called chunking. The span annotation element itself is chunk.

```

<s xml:id="example.p.1.s.1">
  <t>The Dalai Lama greeted him.</t>
  <w xml:id="example.p.1.s.1.w.1"><t>The</t></w>
  <w xml:id="example.p.1.s.1.w.2"><t>Dalai</t></w>
  <w xml:id="example.p.1.s.1.w.3"><t>Lama</t></w>
  <w xml:id="example.p.1.s.1.w.4"><t>greeted</t></w>
  <w xml:id="example.p.1.s.1.w.5"><t>him</t></w>
  <w xml:id="example.p.1.s.1.w.6"><t>.</t></w>
  <chunking>
    <chunk xml:id="example.p.1.s.1.chunk.1">
      <wref id="example.p.1.s.1.w.1" t="The" />
      <wref id="example.p.1.s.1.w.2" t="Dalai" />
      <wref id="example.p.1.s.1.w.3" t="Lama" />
    </chunk>
    <chunk xml:id="example.p.1.s.1.chunk.2">
      <wref id="example.p.1.s.1.w.4" t="greeted" />
    </chunk>
    <chunk xml:id="example.p.1.s.1.chunk.3">
      <wref id="example.p.1.s.1.w.5" t="him" />
      <wref id="example.p.1.s.1.w.6" t="." />
    </chunk>
  </chunking>
</s>

```

The declaration:

```

<annotations>
  <chunking-annotation set="http://ilk.uvt.nl/foia/sets/syntax-nl" />
</annotations>

```

2.9.5 Events as span annotation

We already saw events as structure annotation in Section 2.6.4. But for more fine grained control of timing information a span annotation element in an offset layer is more suited. FoLiA introduces a timing layer in which `timedevent` can be used to span over tokens that form a particular event. Consider the following example:

```
<s>
  <w xml:id="example.p.1.s.1.w.1"><t>I</t></w>
  <w xml:id="example.p.1.s.1.w.2"><t>think</t></w>
  <w xml:id="example.p.1.s.1.w.3"><t>I</t></w>
  <w xml:id="example.p.1.s.1.w.4"><t>have</t></w>
  <w xml:id="example.p.1.s.1.w.5"><t>to</t></w>
  <w xml:id="example.p.1.s.1.w.6"><t>go</t></w>
  <w xml:id="example.p.1.s.1.w.7"><t>.</t></w>
  <timing>
    <timedevent class="utterance" begintime="2011-12-15T19:01"
      endtime="2011-12-15T19:03" actor="myself">
      <wref id="example.p.1.s.1.w.1" t="I" />
      <wref id="example.p.1.s.1.w.2" t="think" />
    </timedevent>
    <timedevent class="cough" begintime="2011-12-15T19:03"
      endtime="2011-12-15T19:05" actor="myself">
    </timedevent>
    <timedevent class="utterance" begintime="2011-12-15T19:05"
      endtime="2011-12-15T19:06" actor="myself">
      <wref id="example.p.1.s.1.w.3" t="I" />
      <wref id="example.p.1.s.1.w.4" t="have" />
      <wref id="example.p.1.s.1.w.5" t="to" />
      <wref id="example.p.1.s.1.w.6" t="go" />
    </timedevent>
  </timing>
</s>
```

These timed events may also be nested. As always, the classes in the example are set-defined rather than predefined by FoLiA. The (optional) features `begintime` and `endtime` can be used express the exact moment at which an event started or ended. Note that this differs from the generic `datetime` attribute, which would describe the time at which the annotation was recorded, rather than when the event took place!

The declaration:

```
<annotations>
  <timedevent-annotation set="http://ilk.uvt.nl/folia/sets/events" />
</annotations>
```


For more coarser annotation of events, use the structure annotation element `event` instead. See section 2.6.4.

2.9.6 Semantic roles

Development Notes

Still to be done.. The `semroles` layer and `semrole` span annotation element will be reserved for this.

2.9.7 Alternative Span Annotations

With token annotations one could specify an unbounded number of alternative annotations. This is possible for span annotations as well, but due to the different nature of span annotations this happens in a slightly different way.

Where we used `alt` for token annotations, we now use `altlayers` for span annotations. Under this element several alternative layers can be presented. Analogous to `alt`, any layers grouped together are assumed to be somehow dependent. Multiple `altlayers` can be added to introduce independent alternatives. Each alternative should be associated with a unique identifier, which uses “alt” rather than “altlayers”. In this case, any layers it consists of must be marked as non-authoritative using `auth="no"`.

Below is an example of a sentence that is chunked in two ways:

```
<s xml:id="example.p.1.s.1">
  <t>The Dalai Lama greeted him.</t>
  <w xml:id="example.p.1.s.1.w.1"><t>The</t></w>
  <w xml:id="example.p.1.s.1.w.2"><t>Dalai</t></w>
  <w xml:id="example.p.1.s.1.w.3"><t>Lama</t></w>
  <w xml:id="example.p.1.s.1.w.4"><t>greeted</t></w>
  <w xml:id="example.p.1.s.1.w.5"><t>him</t></w>
  <w xml:id="example.p.1.s.1.w.6"><t>.</t></w>
  <chunking>
    <chunk xml:id="example.p.1.s.1.chunk.1">
      <wref id="example.p.1.s.1.w.1" t="The" />
      <wref id="example.p.1.s.1.w.2" t="Dalai" />
      <wref id="example.p.1.s.1.w.3" t="Lama" />
    </chunk>
    <chunk xml:id="example.p.1.s.1.chunk.2">
      <wref id="example.p.1.s.1.w.4" t="greeted" />
    </chunk>
  </chunking>
</s>
```

```

    <chunk xml:id="example.p.1.s.1.chunk.3">
      <wref id="example.p.1.s.1.w.5" t="him" />
      <wref id="example.p.1.s.1.w.6" t="." />
    </chunk>
  </chunking>
  <altlayers xml:id="example.p.1.s.1.alt.1">
    <chunking annotator="John□Doe"
      annotortype="manual" confidence="0.0001" auth="no">
      <chunk xml:id="example.p.1.s.1.alt.1.chunk.1">
        <wref id="example.p.1.s.1.w.1" t="The" />
        <wref id="example.p.1.s.1.w.2" t="Dalai" />
      </chunk>
      <chunk xml:id="example.p.1.s.1.alt.1.chunk.2">
        <wref id="example.p.1.s.1.w.2" t="Lama" />
        <wref id="example.p.1.s.1.w.4" t="greeted" />
      </chunk>
      <chunk xml:id="example.p.1.s.1.alt.1.chunk.3">
        <wref id="example.p.1.s.1.w.5" t="him" />
        <wref id="example.p.1.s.1.w.6" t="." />
      </chunk>
    </chunking>
  </altlayers>
</s>

```

The support for alternatives and the fact that multiple layers (including those of different types) cannot be nested in a single inline structure, should make clear why FoLiA uses a stand-off notation alongside an inline notation.

2.10 Advanced Paradigm

We introduced the FoLiA paradigm in section 2.4. Now we will introduce some of the more advanced aspects of the FoLiA paradigm. These are relevant especially if you want to submit suggestions for extending the FoLiA format with annotations not yet supported.

2.10.1 Human readable Descriptions

Any token annotation element or span annotation element may hold a desc element containing a human readable description for the annotation. An example of this has been already shown for the sense and gap elements.

2.10.2 Text content and multiple classes

In section 2.3 we have seen the text content element `t`. This element can be associated with structural elements such as `w`, `s`, and `p`. The `offset` attribute may be used to explicitly link the text between child and parent. This is demonstrated on three levels in the following example:

```
<p xml:id="example.p.1">
  <t>Hello. This is a sentence. Bye!</t>
  <s xml:id="example.p.1.s.1">
    <t offset="7">This is a sentence.</t>
    <w xml:id="example.p.1.s.1.w.1"><t offset="0">This</t></w>
    <w xml:id="example.p.1.s.1.w.2"><t offset="5">is</t></w>
    <w xml:id="example.p.1.s.1.w.3"><t offset="8">a</t></w>
    <w xml:id="example.p.1.s.1.w.4" space="no"><t offset="10">sentence</t></w>
    <w xml:id="example.p.1.s.1.w.5"><t offset="18">.</t></w>
  </s>
</p>
```

Moreover, we have seen the `space` attribute, which is a simple alternative that can be used to reconstruct the untokenised text if it is not explicitly provided in a parent's `t` element. Allowed values for `space` are:

- “yes” or “ ” (a space) – This is the default and says that the token is followed by a single space.
- “no” or “” (empty) – This states that the token is not followed by a space.
- any other character or string – This states that the token is followed by another character or string that acts as a token separator.

When explicit text content on sentence/paragraph level is provided, offsets can be used to refer back to it from deeper text-content elements. This does imply that there are some challenges to solve. First of all, by default, the `offset` refers to the direct parent of whatever element the text content (`t`) is a member of. If a level is missing we have to explicitly specify this reference using the `ref` attribute. Note that there is no text content for the sentence in the following example, and we refer directly to the paragraph's text:

```
<p xml:id="example.p.1">
  <t>Hello. This is a sentence. Bye!</t>
  <s xml:id="example.p.1.s.1">
    <w xml:id="example.p.1.s.1.w.1">
```

```

        <t ref="example.p.1" offset="7">This</t>
    </w>
    <w xml:id="example.p.1.s.1.w.2">
        <t ref="example.p.1" offset="12">is</t>
    </w>
    <w xml:id="example.p.1.s.1.w.3">
        <t ref="example.p.1" offset="15">a</t>
    </w>
    <w xml:id="example.p.1.s.1.w.4" space="no">
        <t ref="example.p.1" offset="17">sentence</t>
    </w>
    <w xml:id="example.p.1.s.1.w.5">
        <t ref="example.p.1" offset="25">.</t>
    </w>
</s>
</p>

```

Note that text content is always expected to be untokenised, except in w tags as it by definition is the tokenisation layer.

It is possible to associate *multiple text-content* elements with the same element, and thus associating multiple texts with the same element. You may wonder what could possibly be the point of such extra complexity. But there is a clear use-case when dealing with for example corrections, or wanting to associate the text version just prior or after a processing step such as Optical Character Recognition or another kind of normalisation.

Corrections are challenging because they can be applied to text content and thus change the text. Corrections are often applied on the token level (within w tags), but you may want them propagated to the text content of sentences or paragraphs whilst at the same time wanting to retain the text how it originally was. This can be accomplished by introducing text content of a different class. Text content that has no associated class obtains the “current” class by default, it is expected to always be up-to-date. There is a notable exception: text content that appears within the scope of original elements within a correction element automatically adopts the “original” class.² This thus implies that in this rare case, FoLiA actually pre-defines classes (i.e: “original” and “current”)! In addition to these two pre-defined classes, any other custom classes may be added as you see fit. If you add custom classes, you need a declaration:

```

<annotations>
    <text-annotation set="http://ilk.uvt.nl/folia/sets/text-annotation" />

```

²For more deeply nested original elements, you will have to assign your own classes if you do not want them to take the “original” class.

</annotations>

Below is an example illustrating the usage of multiple classes. To show the flexibility, offsets are added, but these are of course always optional. Note that when an offset is specified, it always refers to a text-content element of the same class!

```
<p xml:id="example.p.1">
  <t>Hello. This is a sentence. Bye!</t>
  <t class="original">Hello. This iz a sentence. Bye!</t>
  <s xml:id="example.p.1.s.1">
    <t offset="7">This is a sentence.</t>
    <t class="original" offset="7">This is a sentence.</t>
    <w xml:id="example.p.1.s.1.w.1">
      <t offset="0">This</t>
    </w>
    <w xml:id="example.p.1.s.1.w.2">
      <correction>
        <new>
          <t offset="5">is</t>
        </new>
        <original auth="no">
          <t offset="5">iz</t>
          <!-- Note that this element has class 'original' by definition! -->
        </original>
      </correction>
    </w>
    <w xml:id="example.p.1.s.1.w.3">
      <t offset="8">a</t>
    </w>
    <w xml:id="example.p.1.s.1.w.4" space="no">
      <t offset="10">sentence</t>
    </w>
    <w xml:id="example.p.1.s.1.w.5">
      <t offset="48">.</t>
    </w>
  </s>
</p>
```

One important aspect is that FoLiA dictates that `t` elements without classes always appear first, before any other text-content tags that do specify classes. This rule facilitates the job of parsers in quickly getting the latest text content even if there are multiple classes: the first `t` element will be the most recent by definition.

In the above example, the correction is explicit, in the next example, it is implicit. Furthermore, to illustrate how you could use other custom classes, the next

example introduces an custom “ocroutput” class that shows the (fictious) output of an OCR system prior to some implicit correction stage.

```
<p xml:id="example.p.1">
  <t>Hello. This is a sentence. Bye!</t>
  <t class="original">Hello. This iz a sentence. Bye!</t>
  <t class="ocroutput">Hell0 Th1s iz a sentence, Bye1</t>
  <s xml:id="example.p.1.s.1">
    <t offset="7">This is a sentence.</t>
    <t class="original" offset="7">This is a sentence.</t>
    <t class="ocroutput" offset="6">Th1s iz a sentence,</t>
    <w xml:id="example.p.1.s.1.w.1">
      <t offset="0">This</t>
      <t class="ocroutput" offset="0">Th1s</t>
    </w>
    <w xml:id="example.p.1.s.1.w.2">
      <t offset="5">is</t>
      <t offset="5" class="original">iz</t>
      <t offset="5" class="ocroutput">iz</t>
    </w>
    <w xml:id="example.p.1.s.1.w.3">
      <t offset="8">a</t>
    </w>
    <w xml:id="example.p.1.s.1.w.4" space="no">
      <t offset="10">sentence</t>
    </w>
    <w xml:id="example.p.1.s.1.w.5">
      <t offset="48">.</t>
      <t offset="48" class="original">.</t>
      <t offset="48" class="ocroutput">,</t>
    </w>
  </s>
</p>
```

Last, an important note regarding offsets; all offset values are measured in uni-code code-points, the first character having index zero. Take special care with combining diacritical marks versus codepoints that directly integrate the diacritical mark.

2.10.3 Features, Subsets and multiple classes

For some annotations, associating a single class from a set is not sufficient. FoLiA provides the option to associate multiple classes, or to associate classes with specific subsets.

In addition to a main class, an arbitrary number of *features* can be added to *any* annotation element. Each feature pertains to a specific *subset*. Subsets and the classes within them can be invented at will as they are part of the set definition, which is left entirely to the user. However, certain annotation elements also have some predefined subsets you may use.

The element `feat` is used to add features to any kind of annotation. In the following example we make use of a subset we invented which ties a lemma to a page number in some dictionary where the lemma can be found.

```
<lemma class="house">
  <feat subset="dictionary_page" class="45" />
</lemma>
```

A more thorough example for Part-of-Speech tags with features will be explained in section .

Some elements have predefined subsets because some features are very commonly used. It however still depends on the set on whether these can be used, and which values these take. Whenever subsets are predefined they can be assigned using XML attributes. Consider the following example of lexical semantic sense annotation, in which subset “synset” is predefined:

```
<sense class="X" synset="Y" />
```

This is semantically equivalent to:

```
<sense class="X">
  <feat subset="synset" class="Y" />
</sense>
```

The following example of event annotation with the feature with predefined subset “actor” is similar:

```
<event class="tweet" actor="John_Doe">
  ...
</event>
```

This is semantically equivalent to:

```
<event class="tweet">
  <feat subset="actor" class="John_Doe" />
  <s>...</s>
</event>
```

Features can also be used to assign multiple classes within the same subset, which is impossible with main classes. In the following example the event is associated with a list of two actors. In this case the XML attribute shortcut no longer suffices, and the `feat` element must be used.

```
<event class="conversation">
  <feat subset="actor" class="John_Doe" />
  <feat subset="actor" class="Jane_Doe" />
  <p>...</p>
</event>
```

To recap: the `feat` element can always be used to freely to associate any additional classes of *any* designed subset with with *any* annotation element. For certain elements, there are predefined subsets, in which case you can assign them using the XML attribute shortcut. This however only applies to the pre-defined features.

2.10.4 Part-of-Speech tags with features

Part-of-speech tags are a good example of the scenario outlined above. Part-of-speech tags may consist of multiple features, which in turn *may* be associated with specific subsets. There are two scenarios envisionable, one in which the class of the `pos` element combines all features, and one in which it is the foundation upon which is expanded. Which one is used is entirely up to the defined set.

Option one:

```
<w xml:id="example.p.1.s.1.w.2">
  <t>boot</t>
  <pos head="N" class="N(soort,ev,basis,zijd,stan)">
    <desc>Noun, singular, neuter</desc>
    <feat subset="ntype" class="soort" />
    <feat subset="number" class="ev" />
    <feat subset="degree" class="basis" />
    <feat subset="gender" class="zijd" />
    <feat subset="case" class="stan" />
  </pos>
</w>
```

In FoLiA, this attribute `head` is implicitly associated with the subset “head” of whatever set you defined. This would thus be equal to:

```
<feat subset="head" class="N" />
```


Option two:

```
<w xml:id="example.p.1.s.1.w.2">
  <t>boot</t>
  <pos class="N">
    <desc>Noun, singular, neuter</desc>
    <feat subset="ntype" class="soort" />
    <feat subset="number" class="ev" />
    <feat subset="degree" class="basis" />
    <feat subset="gender" class="zijd" />
    <feat subset="case" class="stan" />
  </pos>
</w>
```

2.11 Subtoken Annotation

Subtoken category combines aspects from both token annotation and span annotation. It introduces an annotation layer at the token level, as child of the `w` element. Within this annotation layer, subtoken annotation elements can be used to annotate parts within the token itself. Whereas span elements use `wref` to refer to words, subannotation elements use `t` to refer to part of the text of the token. Recall that `t` elements can contain references to higher-level `t` elements. In such cases, the `offset` attribute is used to designate the offset index in the word's associated text element (`t`) (zero being right at the start of the text).

2.11.1 Morphological Analysis

```
<w xml:id="example.p.4.s.2.w.4">
  <t>leest</t>
  <lemma class="lezen" />
  <morphology>
    <morpheme>
      <feat subset="type" class="stem">
        <feat subset="function" class="lexical">
          <t offset="0">lees</t>
        </morpheme>
      </morpheme>
      <feat subset="type" class="suffix">
        <feat subset="function" class="inflexional">
          <t offset="4">t</t>
        </morpheme>
      </morpheme>
    </morphology>
  </w>
```

```
</w>
```

```
<annotations>  
  <morphology-annotation set="http://ilk.uvt.nl/folia/sets/entities" />  
</annotations>
```

2.11.2 Named Entities within a token

Named entities may sometimes occur as *part* of a token. The subentities annotation layer and the subentity subtoken annotation element can be used for annotating these:

```
<w xml:id="example.p.4.s.2.w.4">  
  <t>CDA-voorzitter</t>  
  <subentities>  
    <subentity class="org" annotator="Maarten_van_Gompel"  
      annotortype="manual">  
      <t offset="0">CDA</t>  
    </subentity>  
  </subentities>  
</w>
```

The declaration:

```
<annotations>  
  <subentity-annotation set="http://ilk.uvt.nl/folia/sets/entities" />  
</annotations>
```

2.12 Alignments

FoLiA provides a facility to align parts of your document with other parts of your document, or even with parts of other FoLiA documents. These are called *alignments* and are implemented using either the `alignment` element, a form of extended token annotation (that may thus also be applied on other levels such as sentences or paragraphs), or `complexalignment`, a form of span annotation for more complex word alignments.

Consider the following two aligned sentences from two *distinct* FoLiA documents in different languages. Note that the `t` attribute to the `aref` element is merely optional and is this overhead is added simply to facilitate the job of limited FoLiA parsers and providing a quick reference to the target text for both parser and

human user. The `xlink:href` attribute is used to link to the target document, if any. If the alignment is within the same document then it should be simply omitted. The `type` attribute is mandatory and specifies the type of element the alignment points too, i.e. its value is equal to the tagname it points to.

```
<s xml:id="example-english.p.1.s.1">
  <t>The Dalai Lama greeted him.</t>
  <alignment class="french-translation" xlink:href="doc-french.xml"
    xlink:type="simple">
    <aref id="doc-french.p.1.s.1" t="Le_Dalai_Lama_le_saluait." type="s" />
  </alignment>
</s>

<s xml:id="example-french.p.1.s.1">
  <t>Le Dalai Lama le saluait.</t>
  <alignment class="english-translation" xlink:href="doc-english.xml"
    xlink:type="simple">
    <aref id="doc-english.p.1.s.1" t="The_Dalai_Lama_greeted_him." type="s" />
  </alignment>
  <alignment class="english-translation" xlink:href="doc-dutch.xml"
    xlink:type="simple">
    <aref id="doc-dutch.p.1.s.1" t="De_Dalai_Lama_begroette_hem." type="s" />
  </alignment>
</s>
```

Although the above example has a single alignment reference (`aref`), it is not forbidden to link to specify multiple references within the alignment block. For more complex alignments however, such as word alignments that include many-to-one, one-to-many or many-to-many alignments, the `span` annotation element `complexalignment` may be more suitable. The following example illustrates a many-to-many word-alignment of the word “Dalai Lama” from English to French. This takes places within the `complexalignments` annotation layer.

```
<s xml:id="example-english.p.1.s.1">
  <t>The Dalai Lama greeted him.</t>
  <w xml:id="example-english.p.1.s.1.w.1"><t>The</t></w>
  <w xml:id="example-english.p.1.s.1.w.2"><t>Dalai</t></w>
  <w xml:id="example-english.p.1.s.1.w.3"><t>Lama</t></w>
  <w xml:id="example-english.p.1.s.1.w.4"><t>greeted</t></w>
  <w xml:id="example-english.p.1.s.1.w.5"><t>him</t></w>
  <w xml:id="example-english.p.1.s.1.w.6"><t>.</t></w>
  <complexalignments>
    ...
    <complexalignment class="french-translation" xlink:href="doc-french.xml"
      xlink:type="simple">
      <wref id="example-english.p.1.s.1.w.2" t="Dalai">
      <wref id="example-english.p.1.s.1.w.3" t="Lama">
      <aref id="example-french.p.1.s.1.w.2" t="Dalai" type="w">
```

```

        <aref id="example-french.p.1.s.1.w.3" t="Lama" type="w">
    </complexalignment>
    ...
</complexalignments>
</s>

```

2.13 Metadata

FoLiA has support for metadata, most notably the extensive and mandatory declaration section for all used annotations which you have seen throughout this documentation. To complement this, there is FoLiA's native metadata system, in which simple metadata fields can be defined and used at will. FoLiA is also able to operate with IMDI or CMDI metadata, either in external file or stored inline. Note however that storing CMDI or IMDI inside your FoLiA document may cause problems when you want to validate your FoLiA document. It is also incompatible with CMDI or IMDI editors that are unaware of FoLiA.

Reference to CMDI in external file proceeds in the following simple fashion:

```

<metadata type="cmdi" src="/path/or/url/to/metadata.cmdi">
    ...
</metadata>

```

The procedure for IMDI is the same:

```

<metadata type="imdi" src="/path/or/url/to/metadata.imdi">
    ...
</metadata>

```

If you use neither CMDI nor IMDI, then you can use FoLiA's native system, which is very simple: It introduces the `meta` element that allows you to define key value pairs as follows:

```

<metadata type="native">
    <annotations>
        ..
    </annotations>
    <meta id="title">Title to my document</meta>
    <meta id="language">eng</meta>
</metadata>

```

You can simply define fields with custom IDs, but the following fields are pre-defined and recommended to be filled:

- **title** – The title of the FoLiA document
- **language** – An ISO-639-3 language code identifying the language the document is
- **date** – The date of publication in YYYY-MM-DD format
- **publisher** – The publishing institution or individual
- **license** – The type of license of the document (for example: *GNU Free Documentation License*)

Chapter 3

Set Definition Format

3.1 Introduction

The FoLiA format consists not just out of the Document Format discussed in the previous chapter, but also of a Set Definition Format. The document format is agnostic about all sets and the classes therein, it is the Set Definition Format that defines precisely what classes are allowed in a certain set, including any subsets.

Recall from section ?? that all used sets need to be declared in the document header and that they point to URLs holding a FoLiA set definition. If no set definition files are associated, then a full in-depth validation can not take place.

3.2 Types and classes

The set definition format is fairly straightforward, each set definition file represents one set, including all of its subsets.

Here is a simple example:

```
<set xml:id="simplepos" type="closed">
  <class xml:id="N" label="Noun" />
  <class xml:id="V" label="Verb" />
  <class xml:id="A" label="Adjective" />
</set>
```

The ID of the class determines a value the `class` attribute may take in the FoLiA document, for elements of this set. The `label` attribute carries a human readable description for presentational purposes, this is optional but highly recommended.

There are three possible types for sets and subsets:

1. **open**: classes may be anything and are not defined
2. **closed**: classes are defined strictly
3. **mixed**: classes may be anything, but some are predefined

A set definition file for an open type set definition may be as concise as:

```
<set xml:id="lemmas-nl" type="open" />
```

3.3 Concept link

You may want to associate classes, or even sets themselves, with some kind of semantic web or category registry. This link can be made using the `conceptlink` attribute which may be placed on classes, sets and subset elements.

```
<set xml:id="simplepos" type="closed" conceptlink="http://some/host/simplepos">
  <class xml:id="N" label="Noun" conceptlink="http://some/host/noun" />
  <class xml:id="V" label="Verb" conceptlink="http://some/host/verb" />
  <class xml:id="A" label="Adjective" conceptlink="http://some/host/adj" />
</set>
```

FoLiA does not dictate any format requirements for conceptual links, it can be anything, such as an RDF resource, or any other kind. If you want something more specific, or you want to link to multiple semantic resources, simply use your own “`conceptlink`” attribute in a different custom XML namespace.

3.4 Subsets

Section 2.10.3 introduced subsets. These can be defined in a similar fashion to sets and also carry a `type` attribute:

```

<set xml:id="simplepos" type="closed">
  <class xml:id="N" label="Noun" />
  <class xml:id="V" label="Verb" />
  <class xml:id="A" label="Adjective" />
  <subset xml:id="gender" class="closed">
    <class xml:id="m" label="Masculine" />
    <class xml:id="f" label="Feminine" />
    <class xml:id="n" label="Neuter" />
  </subset>
</set>

```

It is possible for subsets to be used multiple times if the subset is declared with the attribute `multi` set to `true` (defaults to `false`). This allows multiple classes to be associated with a subset. Subsets can be made mandatory by setting the attribute `required` to `true`.

3.5 Constraints

Not all classes in subsets can be combined with others. Often the need arises to put constraints on which classes can go together. The previous example already illustrates this. For many languages, the “gender” subset does not make sense on verbs. We can put a constraint on the usage of this subset, limiting its usage to nouns and adjectives:

```

<set xml:id="simplepos" type="closed">
  <class xml:id="N" label="Noun" />
  <class xml:id="V" label="Verb" />
  <class xml:id="A" label="Adjective" />
  <subset xml:id="gender" class="closed">
    <constraint>
      <restrict class="N" />
      <restrict class="A" />
    </constraint>
    <class xml:id="m" label="Masculine" />
    <class xml:id="f" label="Feminine" />
    <class xml:id="n" label="Neuter" />
  </subset>
</set>

```

For sake of brevity, constraints can be named and referred to when they are needed multiple times.

```

<set xml:id="simplepos" type="closed">
  <class xml:id="N" label="Noun" />

```



```

<class xml:id="V" label="Verb" />
<class xml:id="A" label="Adjective" />
<subset xml:id="gender" class="closed">
  <constraint name="constraint.1">
    <restrict class="N" />
    <restrict class="A" />
  </constraint>
  <class xml:id="m" label="Masculine" />
  <class xml:id="f" label="Feminine" />
  <class xml:id="n" label="Neuter" />
</subset>
<subset xml:id="case" class="closed">
  <constraint ref="constraint.1" />
  <class xml:id="nom" label="Nominative" />
  <class xml:id="gen" label="Genitive" />
  <class xml:id="dat" label="Dative" />
  <class xml:id="acc" label="Accusative" />
</subset>
</set>

```

Constraints can be used within subsets, but also within classes:

```

<set xml:id="simplepos" type="closed">
  <class xml:id="N" label="Noun" />
  <class xml:id="V" label="Verb" />
  <class xml:id="A" label="Adjective" />
  <subset xml:id="gender" class="closed">
    <class xml:id="m" label="Masculine">
      <constraint name="constraint.1">
        <restrict class="N" />
        <restrict class="A" />
      </constraint>
    </class>
  </subset>
</set>

```

Using the restrict element, you force a certain class from the main set or any subset, thus enumerating all the allowed classes. For example, the following constraint demands masculine or feminine nouns in either nominative or accusative case. All of the restrictions must be satisfied for the constraint to match, restrictions on the same subset (or the main set if no subset is specified) are automatically considered as disjunctions.

```

<constraint>
  <restrict class="N" />
  <restrict subset="gender" class="f" />
  <restrict subset="gender" class="m" />
  <restrict subset="case" class="nom" />

```

```

    <restrict subset="case" class="acc" />
</constraint>

```

Disjunctions, over different subsets, can be made explicitly using the `disjunction` element. The following would be constrained to feminine nouns *or* plural nouns, rather than feminine plural nouns if the disjunction element were not present:

```

<constraint>
  <restrict class="N" />
  <disjunction>
    <restrict subset="gender" class="f" />
    <restrict subset="number" class="plural" />
  </disjunction>
</constraint>

```

You can also opt to specify the “forbidden” classes using `except`. Only if not a single one of the exceptions applies, the constraint is met.

```

<constraint>
  <except class="V" />
  <except class="A" />
  <except subset="gender" class="n" />
  <except subset="case" class="gen" />
  <except subset="case" class="dat" />
</constraint>

```

Restrict and except elements can also be mixed, in which case the constraint matches if all of the restrictions do, and if none of the exceptions do. Moreover, disjunctions can be nested, to form complex constraints.

The *required* attribute can be used on subsets to indicate whether they are mandatory or optional, but a more powerful mechanism is available using constraints and the `require` element (or its complement: `forbid`). The following example adds a constraint on nouns and requires it to have the gender and number subsets specified, whereas for verbs, tense and number are required.

```

<class="N">
  <constraint>
    <require subset="gender" />
    <require subset="number" />
  </constraint>
</class>
<class="V">
  <constraint>
    <require subset="tense" />
    <require subset="number" />
  </constraint>

```

</class>

Appendix A

Common Queries

Considering the fact that FoLiA is an XML-based format, XPath and its derivatives are the designated tools for searching in a FoLiA document.

A very common XPath predicate found in many XPath expressions for FoLiA is `not(ancestor-or-self::*/@auth)`. This exploits the notion of authoritativeness. Certain elements in FoLiA are non-authoritative, which means that they have no direct bearing on the actual state of the element they describe. The most notable elements that are non-authoritative are alternatives, suggestions for correction, and the original part of a correction. The predicate `not(ancestor-or-self::*/@auth)` guarantees that no elements can be selected that occur within the scope of any non-authoritative element. This prevents selecting for example alternative annotations or annotations that were superseded by a correction step. This is in most cases what the user wants and why you will find this predicate appended to almost every XPath expression for FoLiA.

Some common XPath queries are listed below, note that for brevity and readability the namespace prefix is omitted. In actual situations you will have to specify the FoLiA namespace with each element, as XPath unfortunately has no notion of a default namespace.

- XPath query for all paragraphs: `//p[not(ancestor-or-self::*/@auth)]`
- XPath query for all sentences: `//s[not(ancestor-or-self::*/@auth) and not(ancestor::quote)]`
Explanation: When selecting sentences, you often do not want sub-sentences that are part of a quote, since they may overlap with the larger sentence

they form a part of. The `not(ancestor::quote)` predicate guarantees this can not happen.

- XPath query for all words: `//w[not(ancestor-or-self::*/@auth)]`
- XPath query for the text of all words/tokens: `//w//t[not(ancestor-or-self::*/@auth)][1]`
Explanation: The [1] predicate is important here and makes sure to select only the first text content element in case there are multiple. Recall that FoLiA dictates that the classless text element appears before any others, ensuring that this is the proper way to obtain the most recent text of an element. (See Section 2.10.2.
- XPath query for all words with lemma X: `//w[.//lemma[@class="X" and not(ancestor-or-self::*/@auth)]]` *Note: This query assumes there is only one declaration for lemma annotation, and the set has been verified*
- XPath query for all words with PoS-tag A in set S: `//w[.//pos[@set="S" and @class="A" and not(ancestor-or-self::*/@auth)]]` *Note: This query assumes the set attribute was set explicitly, i.e. there are multiple possible sets in the document for this annotation type*
- XPath query for the text of all words with PoS-tag A in set S: `//w[.//pos[@set="S" and @class="A" and not(ancestor-or-self::*/@auth)]]//t[not(ancestor-or-self::*/@auth)]`
Note: The predicate for non-authoritativeness here needs to be applied both to the pos element and the text content element t, otherwise you may accidentally select the text of words which have the desired pos tag only as an alternative.
- XPath query to select all alternative PoS tags for all words: `//w/alt/pos`

Before you release XPath queries on FoLiA documents, make sure to first parse the declarations present in the metadata (the annotations block). Verify that the annotation type with the desired set you are looking for is actually present, otherwise you needn't bother running a query at all. Note that the XPath expression differs based on whether there is only one set defined for the sought annotation type, or if there are multiple. In the former case, you can't use the `@set` attribute to select, and in the latter case, you must.

Appendix B

Validation

Validation proceeds at two levels: shallow validation and deep validation. Shallow validation considers only the structure of the FoLiA document, without validating the sets and classes used. Deep validation checks the sets and classes for their validity using the set definition files.

Shallow validation is performed using a RelaxNG schema, to be found at <https://github.com/proyc>.

You can validate your document using standard XML tools such as `xmllint` or `jing`, the latter is known to produce friendlier error output in case of validation errors.

```
> xmllint -relaxng folia.rng document.xml  
> jing folia.rng document.xml
```

Development Notes

Deep validation is still being worked on and will most likely use Schematron.

Appendix C

Implementations

The following FoLiA implementations exist currently, both follow a highly object-oriented model in which FoLiA XML elements correspond with classes.

1. `pynlpl.formats.folia` - A FoLiA library in Python. Part of the Python Natural Language Processing Library. Documentation can be found at <http://ilk.uvt.nl/olia/>
2. `libolia` - A FoLiA library in *C++*. (Still under heavily development, September 2011)

The following table shows the level of FoLiA support in these libraries:

	PyNLPI	libfolia	
Programming Language	python	C++	
Token Annotation	yes	yes	
Span Annotation	yes	yes	
Subtoken Annotation	yes	yes	
Subset/feature support	yes	yes	
Alternatives	yes	yes	
Correction Annotation	yes	yes	
Alignments	not yet	not yet	
Text offset support	yes	yes	
Native metadata	yes	?	
Query facility (findwords)	partial ¹	partial ²	
IMDI interpretation	partial ³	no	
RelaxNG schema generation	yes ⁴	no	
RelaxNG validation	yes	no	
Set definition support	partial ⁵	no	
Deep validation (Schematron)	not yet	no	
D-Coi read compatibility	partial ⁶	no	

¹for token annotation only

²for token annotation only

³and only for in-document IMDI

⁴The RelaxNG schema is generated dynamically by the library itself

⁵under development

⁶only basic elements, no List, Figure, etc..

Appendix D

Conversions

from	to	status
D-Coi	FoLiA	yes, via XSLT or pynlpl.1.f
FoLiA	D-Coi	yes, via XSLT ²
Frog/Tadpole columned format	FoLiA	yes, using Python script and
FoLiA	Frog/Tadpole columned format,CSV	yes, using Python script and
FoLiA	xhtml	yes, via XSLT ³
Alpino	FoLiA	not yet
FoLiA	Alpino	not yet

¹via library only with basic elements, no List, Figure, etc..

²With inevitable data loss for constructs that can not be represented in D-Coi

³only support for some token annotations

Bibliography

- [1] Eneko Agirre¹, Xabier Artola¹, Arantza Diaz de Ilarraza¹, German Rigau¹, Aitor Soroa¹, and Wauter Bosma. Kyoto annotation format, 2009.
- [2] Wilko Apperloo. XML basisformaat D-Coi: Voorstel XML formaat presentational markup. Technical report, Polderland Language and Speech Technology, 2006.
- [3] Gosse Bouma, Gertjan van Noord, and Rob Malouf. Alpino: Wide-coverage computational analysis of dutch. In Walter Daelemans, Khalil Sima'an, Jorn Veenstra, and Jakub Zavrel, editors, *CLIN*, volume 37 of *Language and Computers - Studies in Practical Linguistics*, pages 45–59. Rodopi, 2000.
- [4] Tim Bray, Jean Paoli Paoli, and C.M. Sperberg-McQueen. Extensible markup language (XML) 1.0. Technical report, W3C, 2 1998.
- [5] Nancy Ide and Jean Véronis. *The Text Encoding Initiative : Background and Context*. Kluwer Academic Publishers, Dordrecht, 1995.
- [6] Nelleke Oostdijk, Martin Reynaert, Paola Monachesi, Gert-Jan Van Noord, Roeland Ordelman, Ineke Schuurman, and Vincent Vandeghinste. From D-Coi to SoNaR: A reference corpus for dutch. In *Proceedings of the Sixth International Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco, 2008.

Appendix E

Comparison

There are numerous annotation formats in the field, so one might wonder why we felt the need to create yet another one. First of all, there are many old and ad-hoc formats in the computational linguistics community, formats that are often specific to one particular application and lack the expressive power for other kind of annotations than initially designed for. Often, such formats are simple text or column-based formats without any kind of formal scheme. FoLiA is intended as a solution for replacing such formats for the storage and exchange of language resources, and as such to facilitate inter-connectivity of NLP tools and sharing of data.

There are also formats which are modern, XML-based and formalised. In the remainder of this appendix we will draw a quick comparison with some notable formats:

E.0.1 TCF v0.4 – Text Corpus Format

The Text Corpus Format[?] is specifically designed as a language resource exchange format for webservices chained in a workflow, which is also one of the intended uses of FoLiA. The format is a modern, XML-based, unicode encoded, and supports a fair subset of the annotation types supported by FoLiA. It is one of the formats closest to FoLiA and therefore most interesting to compare. All annotations in TCF are represented in a stand-off fashion within the same XML-file, although there seems to be a provision for some inline annotations as well. The format is simpler and less verbose than FoLiA, slimmer and processing

efficiency also explicitly being one of the aims of the authors. This however does compromise expressivity and makes it less expressive than FoLiA. For example, there seems to be no support for corrections, alternatives and the encoding of annotators, time of annotation, and confidence level, something that in FoLiA is deeply integrated into the paradigm. TCF is tagset independent and tagsets can be specified by the user, but it does not generalise or formalise the “set/class” paradigm to the extent FoLiA does. The level of generalisation does not go as far as it does in FoLiA.

In annotation formats there is always a balance between expressivity and efficiency/simplicity. FoLiA is clearly on the more expressive side here, offering more features but therefore making higher demands on memory usage, whereas TCF is more on the efficiency side.

E.0.2 KAF – Kyoto Annotation Format

(Yet to write)

E.0.3 LAF – Linguistic Annotation Framework, and GraF

(Yet to write)

E.0.4 TEI P5 – Text Encoding Initiative

(Yet to write)