

# From Project Data to Sustainable Archiving of Linguistic Corpora

Thomas Schmidt (University of Hamburg) & Andreas Witt (University of Tübingen) <http://www.sfb441.uni-tuebingen.de/c2/index-engl.html>

## Project "Sustainability of Linguistic Resources"

### Project facts

- since 2006

A joint initiative by three German research centres:

- SFB 441 'Linguistic Data Structures' (Tübingen)
- SFB 538 'Multilingualism' (Hamburg)
- SFB 632 'Information Structure' (Potsdam/Berlin)
- over 40 individual research projects
- heterogeneous collection of linguistic corpora: written and spoken, many languages, different research paradigms, different annotation categories, different corpus tools and storage formats

### Project aims

- consolidate different approaches to corpus encoding and processing
- pave the way for sustainable archiving of linguistic resources
- develop rules of best practice for sustainable data handling and sharing

### Seven work packages

- Annotation frameworks
- Query tools
- Tools for data access
- Meta-data
- Integration of terminologies
- General rules of best practice
- Legal issues in data sharing and archiving

## TUSNELDA, EXMARALDA and PAULA

```
<figure id="s35b5" entity="belgiji/s35b5.bmp">
  <figtrans>
    <sp who="Obelika">
      <spokenpar>
        Gde da nasdostrok:em belu zastavu ?
        <marked type="deic-loc">Ovdo</marked>
        ja sive puato !
      </spokenpar>
      <keywords>
        <term>open hands <term>
        <term>lightly bent </term>
      </keywords>
      <situation>
        <sp who="Asterika">
          [...]
        </sp>
      </situation>
    </figtrans>
  </figure>
```



### TUSNELDA (Tübingen)

- XML format, data repository, annotation guideline
- predominantly intended for corpora of *written* language
- *hierarchy-centric* conception of data
- *inline* annotation (single document)
- based on suggestions by the *TEI* and *XCES*

Figure on the left: an Asterix cartoon and its annotation in TUSNELDA

```
Sel He he.
Yil Okul kitabidir. O okül/ okul kitabın içinde
TL-Yil school book-PSS3SG-COP DEI school- book-GEN inside-PSS3SG-U-LOC
Yil [eu] It's a textbook. What's in this textbook?
Sel [k] affirmative
Yil [k] for: kitabının
```

A musical score transcript of a Turkish conversation, generated from EXMARALDA data

```
text.xml
<Body>Fürchtet auch nicht ! Die einstige Fußball-Weltmacht zittert vor einem Winzling . Mit seinem Tor zum 1 : 0 ... </Body>

tok.xml
<Marklist type="token" xml:base="tok.xml">
  <feat xlink:href="#xpointer(id("t2")) values="type infStat-#new"/>
  <mark id="t2" xlink:href="#xpointer(string-range(/body,"",0,8))"/>
  <mark id="t3" xlink:href="#xpointer(string-range(/body,"",14,5))"/>

infStat.xml
<Featlist type="information status" xml:base="tok.xml">
  <feat xlink:href="#xpointer(id("t2")) values="type infStat-#new"/>
  <feat xlink:href="#xpointer(id("t3")) values="type infStat-#new"/>

type_infStat.xml
<typelist type="information status">
  <type id="new" descr="The referent is new in the discourse"/>
```

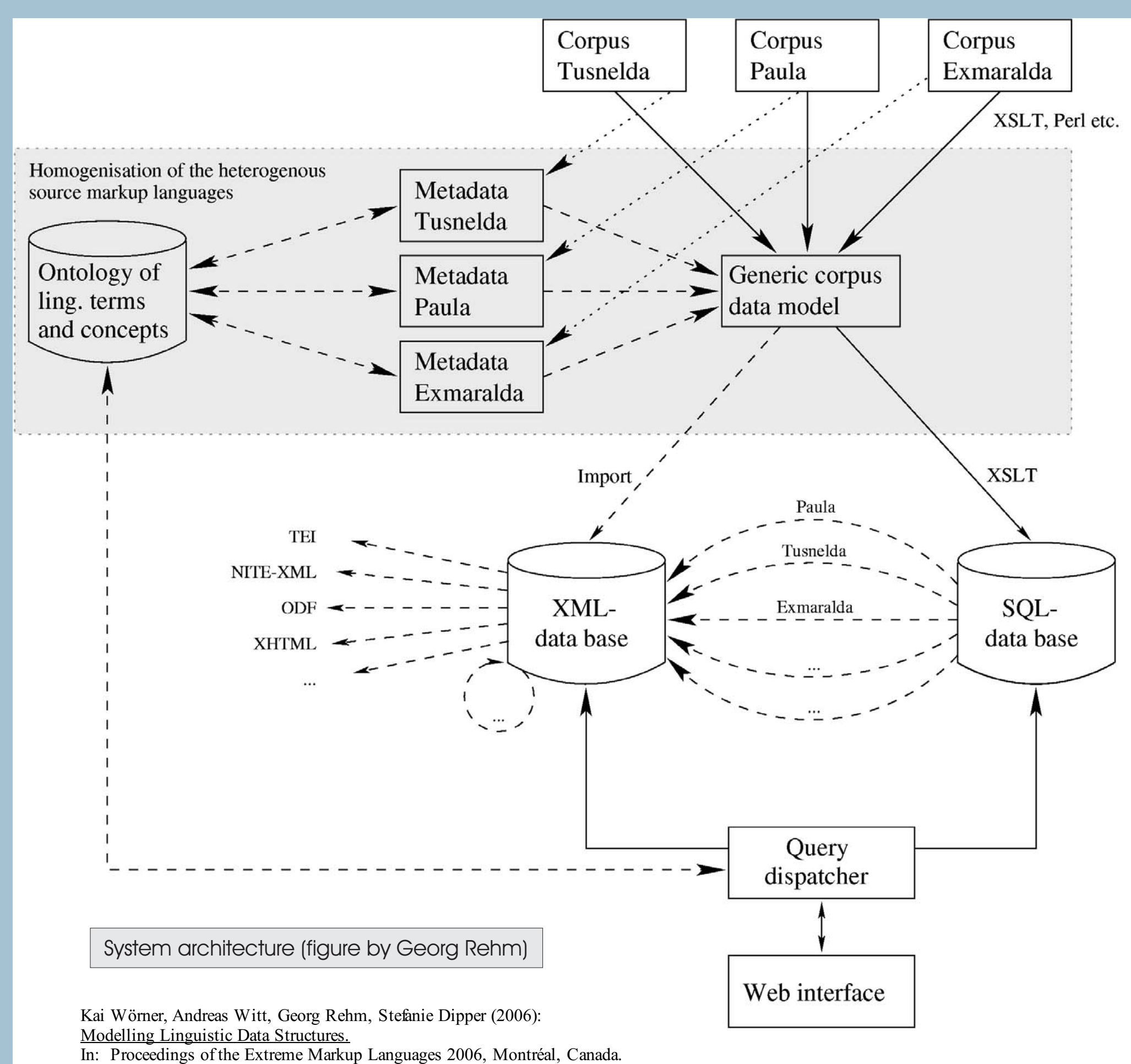
### EXMARALDA (Hamburg)

- XML format, tools for transcription, corpus management, query
- predominantly intended for corpora of *spoken* language
- *timeline-centric* conception of data
- *standoff* annotation (single document)
- based on the idea of *Annotation Graphs*
- EXMARALDA for spoken, *TEI/MENOTA* for written data!

### PAULA (Potsdam/Berlin)

- XML format, data repository, annotation guideline
- for corpora of *spoken and written* language
- hybrid conception of data (multiple hierarchies with timestamps)
- *standoff* annotation (multiple documents)
- based on Linguistic Annotation Framework (LAF)

Figure on the left: an example of PAULA standoff annotation of information structure



## Approach: Generalised data model with different serialisations

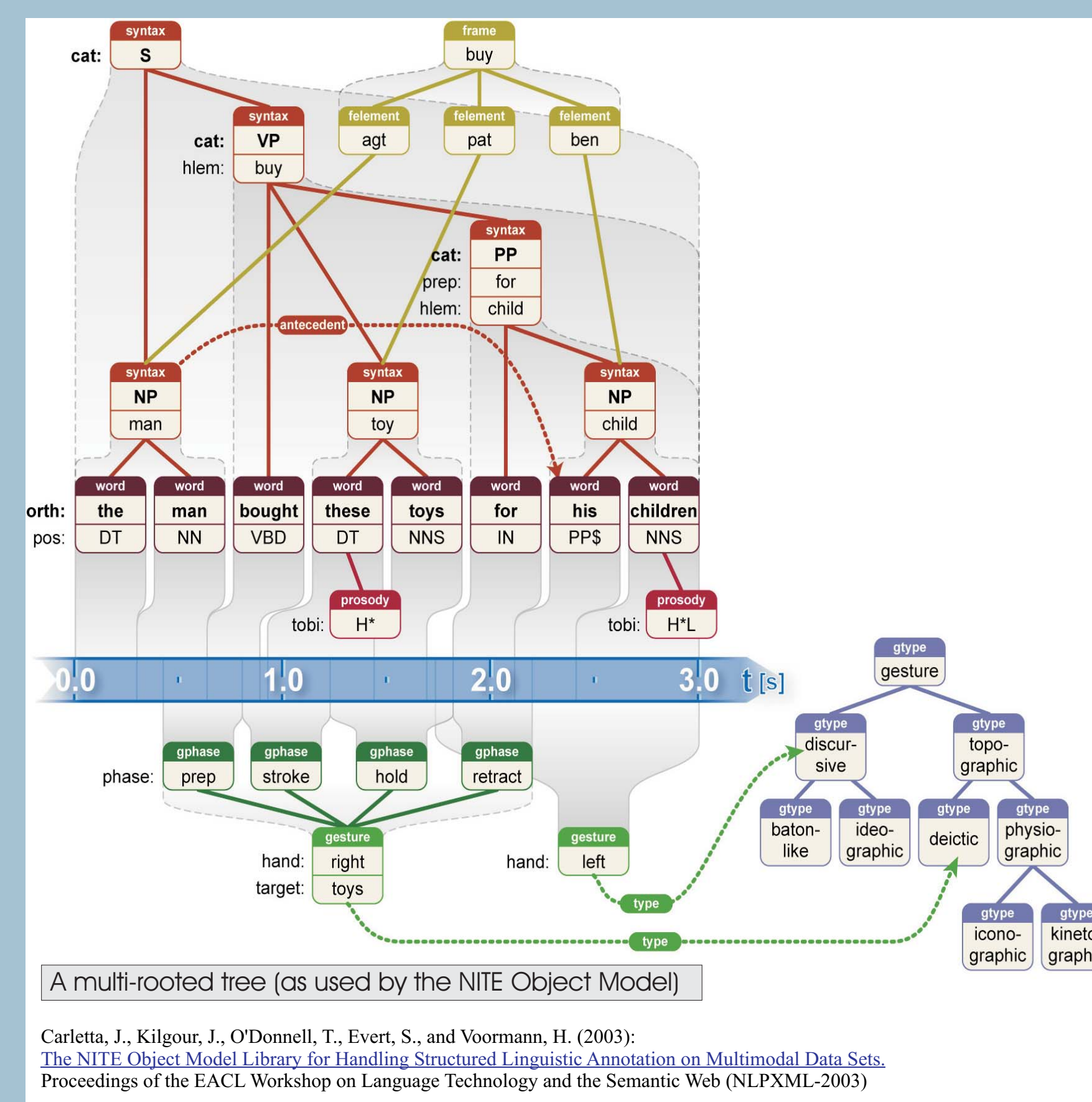
- generalise over existing frameworks
- abstract over physical file formats (serialisations)

### Generic Corpus Data Model (GCDM)

- provide mappings from existing frameworks into GCDM
- provide mappings from GCDM to sustainable serialisations (e.g. TEI conformant XML)

### Properties of the Generic Corpus Data Model

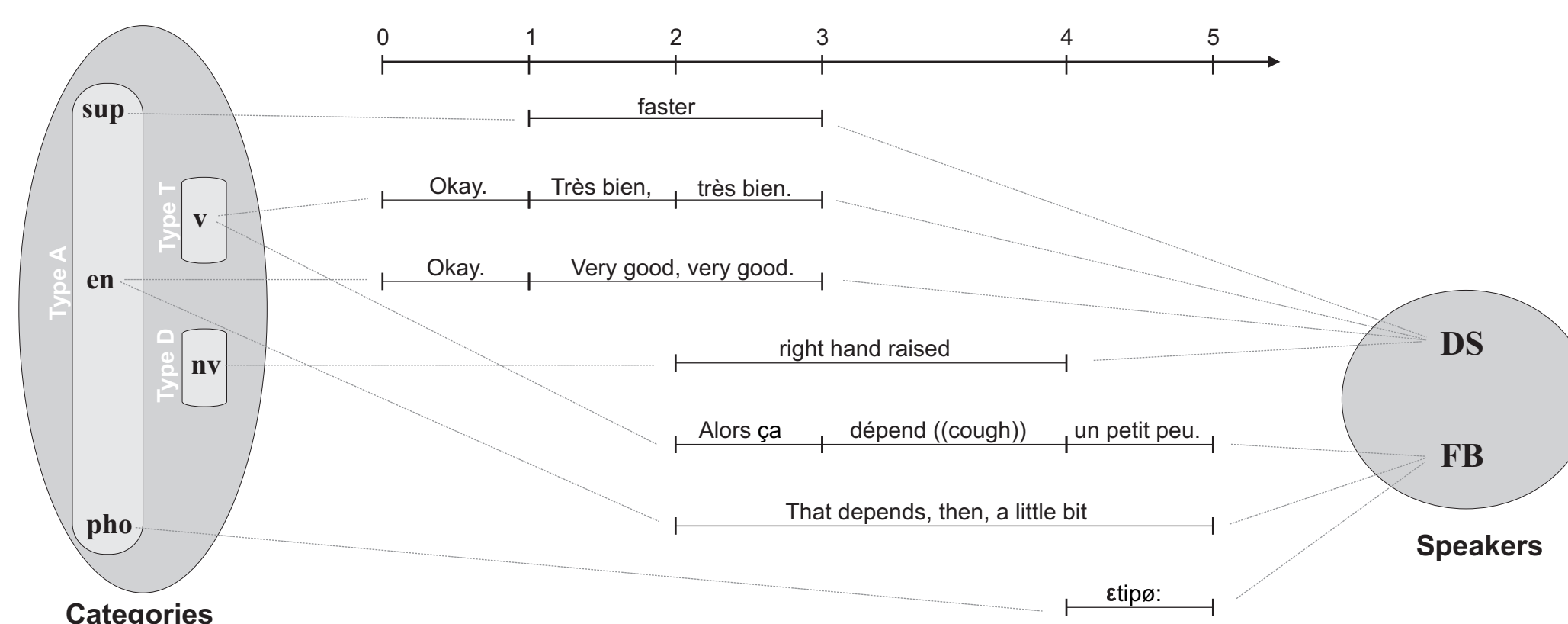
- a multi-rooted tree
- optional time-stamps
- trees relate to another via shared leaves
- similar to the NITE Object Model (Carletta et al. 2003), but: more flexible base level



## Example: Time-based data models and the TEI's Guidelines for Transcriptions of Speech

### The single timeline, multiple tiers data model

- one fully ordered timeline
- events, organised into tiers
- assign each event to a start and end point on the timeline
- no overlap of events within a tier
- assign tiers to speaker, type and category
- used by: PRAAT, EXMARALDA, TASX-Annotator, ELAN, ANVIL, ...



### Serialisation in an EXMARALDA file (stratified)

```
<basic-transcription>
  <head>
    <speakerstale>
      <speaker id="SPK0" abbreviation="DS"/>
      <speaker id="SPK1" abbreviation="FB"/>
    </speakerstale>
  </head>
  <body>
    <common-timeline>
      <tl id="T1"/>
      <tl id="T2"/>
      <!-- [...] -->
      <tl id="T7"/>
    </common-timeline>
    <tier id="T12" speaker="SPK0" category="sup" type="a">
      <event start="T2" end="T4">faster</event>
    </tier>
    <tier id="T122" speaker="SPK0" category="v" type="t">
      <event start="T1" end="T2">Okay. </event>
      <event start="T2" end="T3">Très bien. </event>
      <event start="T3" end="T4">Très bien. </event>
    </tier>
    <tier id="T123" speaker="SPK0" category="en" type="a">
      <event start="T1" end="T2">Okay. </event>
      <event start="T2" end="T4">Very good, very good. </event>
    </tier>
    <tier id="T124" speaker="SPK0" category="nv" type="d">
      <event start="T3" end="T6">right hand raised</event>
    </tier>
    <tier id="T125" speaker="SPK1" category="v" type="t">
      <event start="T3" end="T4">Alors ça </event>
      <event start="T4" end="T5">dépend ((cough)) </event>
      <event start="T5" end="T6">un petit peu. </event>
    </tier>
  </body>
</basic-transcription>
```

Compatibility between the TEI's guidelines for transcriptions of speech and timeline-centric data formats (e.g. Praat, EXMARALDA, ELAN) is achieved via the reference to a **common underlying data model** - the "single timeline, multiple tiers" data model.

Going from the EXMARALDA format to a TEI conformant format is a **hierarchisation**:

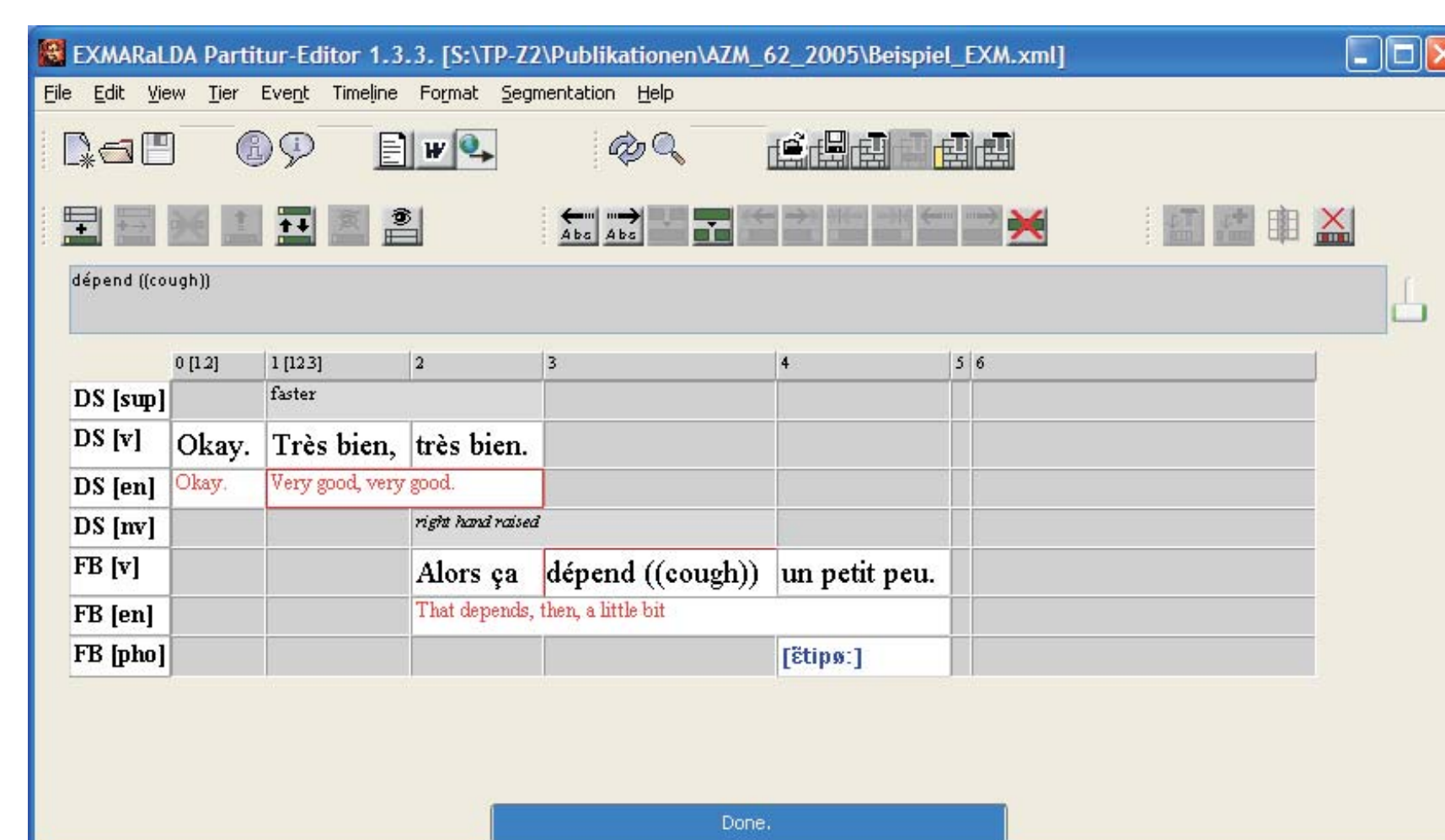
- summarise contiguous sequences of events in tiers of type 'Y' into <u> elements
- use **start** and **end** attributes to assign <u> elements to the timeline
- use <anchor> elements for timeline assignment inside <u> elements
- map other elements accordingly

Going the other way is a **stratification**:

- decompose <u> elements into <event> elements (according to timestamps)
- distribute events onto tiers (one tier per speaker/category combination)

Benefits:

- clearly defined subset of TEI tags for transcriptions of speech
- clearly defined usage of these tags
- use existing transcription software to create or process TEI-conformant data



### Serialisation in a TEI conformant file (hierarchised)

```
<TEI.2>
  <teiHeader>
    <fileDesc/>
    <profileDesc/>
    <particDesc/>
    <person id="DS"/>
    <person id="FB"/>
  </profileDesc>
  </teiHeader>
  <text>
    <timeline>
      <when id="T0"/>
      <when id="T1"/>
      <!-- [...] -->
      <when id="T5"/>
    </timeline>
    <u who="DS" start="T0" end="T3">
      <div type="segmental">
        Okay. <anchor synch="T1"/>
        Très bien. <anchor synch="T2"/>
        très bien.
      </div>
      <div type="prosody">
        <prosody feature="tempo" desc="getting faster" start="T1" end="T3"/>
      </div>
    </u>
    <event who="DS" desc="right hand raised" start="T2" end="T4"/>
    <u who="FB" start="T2" end="T5">
      <div type="segmental">
        Alors ça <anchor synch="T3"/>
        dépend <vocal desc="cough"/> <anchor synch="T4"/>
        un petit peu.
      </div>
    </u>
  </text>
</TEI.2>
```