# Ontology-Based XQuery'ing of XML-Encoded Language Resources on Multiple Annotation Layers

Georg Rehm[1], Richard Eckart[2], Christian Chiarcos[3], Johannes Dellert[1]

University of Tübingen[1]
SFB 441: Linguistic Data Structures
Tübingen, Germany

TU Darmstadt[2]
Dept. of English Linguistics
Darmstadt, Germany

University of Potsdam[3]
SFB 632: Information Structure
Potsdam, Germany

# Context

- Long-term availability of linguistic resources

- Joint Project "Sustainability of Linguistic Data"

- Consolidation of the corpora and data formats

  - Tusnelda       SFB 441  "Linguistic Data Structures"
  - Exmaralda    SFB 538  "Multilingualism"
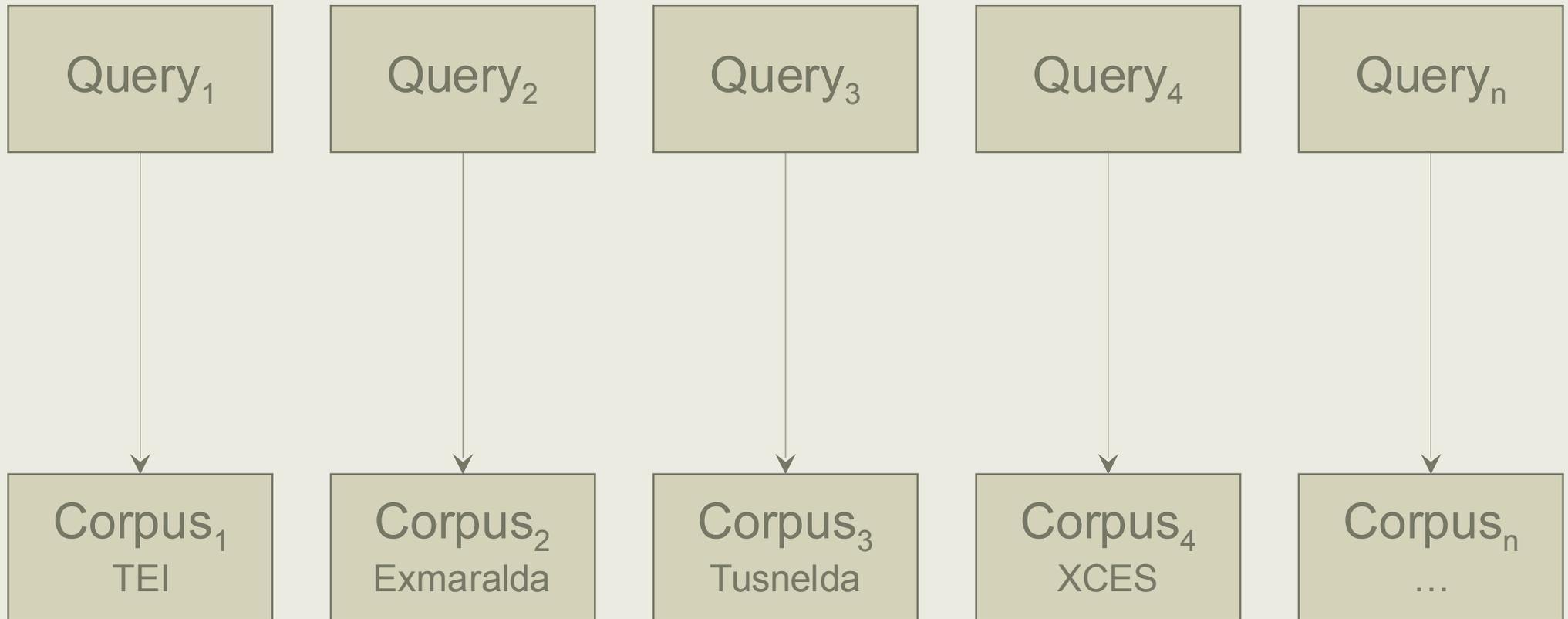  - Paula             SFB 632  "Information Structure"

# SPLICR

- Sustainability Platform for Linguistic Corpora and Resources

    - ~60 highly heterogeneous linguistic resources

- Goals

    - Centralized corpus platform

    - Homogeneous means of accessing and querying

    - Generalisation over

        - Format (Tusnelda, Exmaralda, etc.)

        - Semantics (various tag-sets)

    - Web-based user interface

        - Intuitively usable for linguists
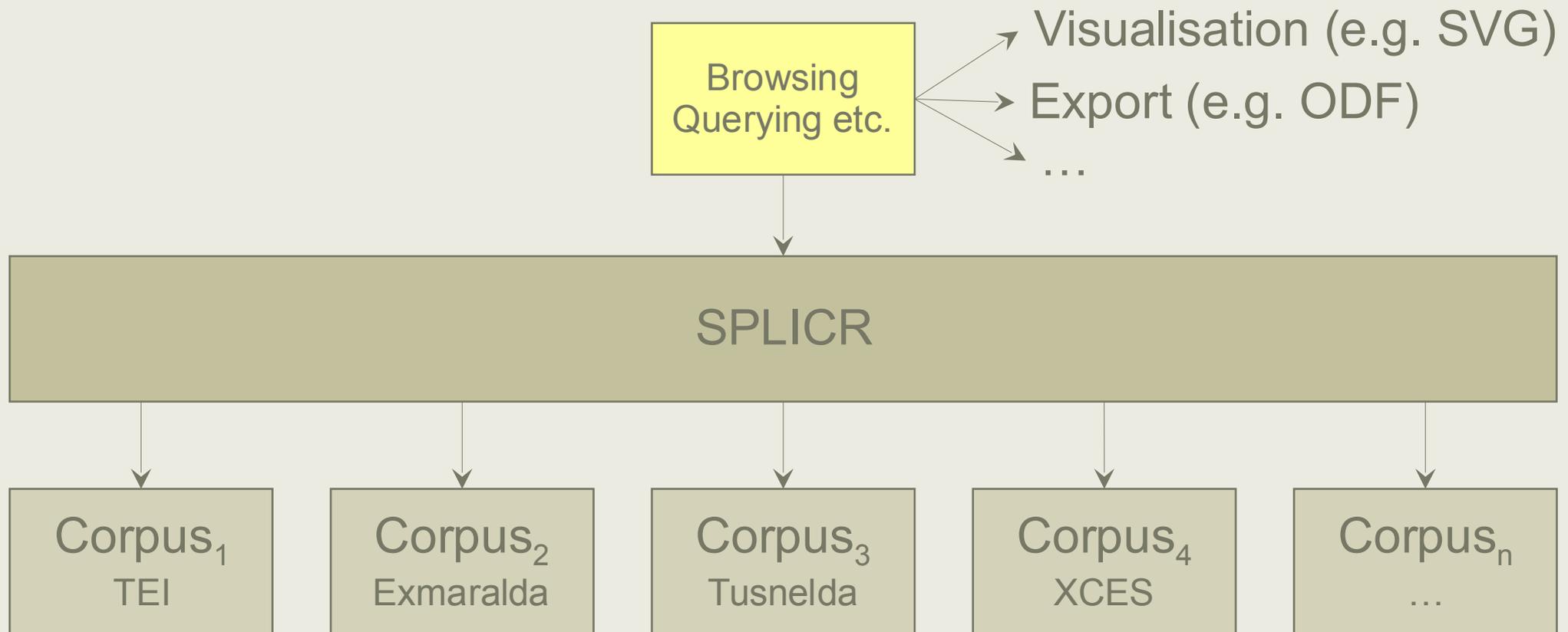
# Linguistic Corpora

- Corpus specific queries

| Query$_1$ | Query$_2$ | Query$_3$ | Query$_4$ | Query$_n$ |

| Corpus$_1$ TEI | Corpus$_2$ Exmaralda | Corpus$_3$ Tusnelda | Corpus$_4$ XCES | Corpus$_n$ ... |

# Linguistic Corpora

- Query against SPLICR

- SPLICR generalises over corpora

- Common visualisation/export modules

*best case scenario*



```
┌─────────────┐
│  Browsing   │ ──→ Visualisation (e.g. SVG)
│ Querying etc.│ ──→ Export (e.g. ODF)
└─────────────┘ ──→ …
       │
       ▼
┌──────────────────────────────────────┐
│               SPLICR                  │
└──────────────────────────────────────┘
```

| Corpus$_1$ | Corpus$_2$ | Corpus$_3$ | Corpus$_4$ | Corpus$_n$ |
|:---:|:---:|:---:|:---:|:---:|
| TEI | Exmaralda | Tusnelda | XCES | … |

Ontology-Based XQuery'ing of XML-Encoded Language Resources on Multiple Annotation Layers
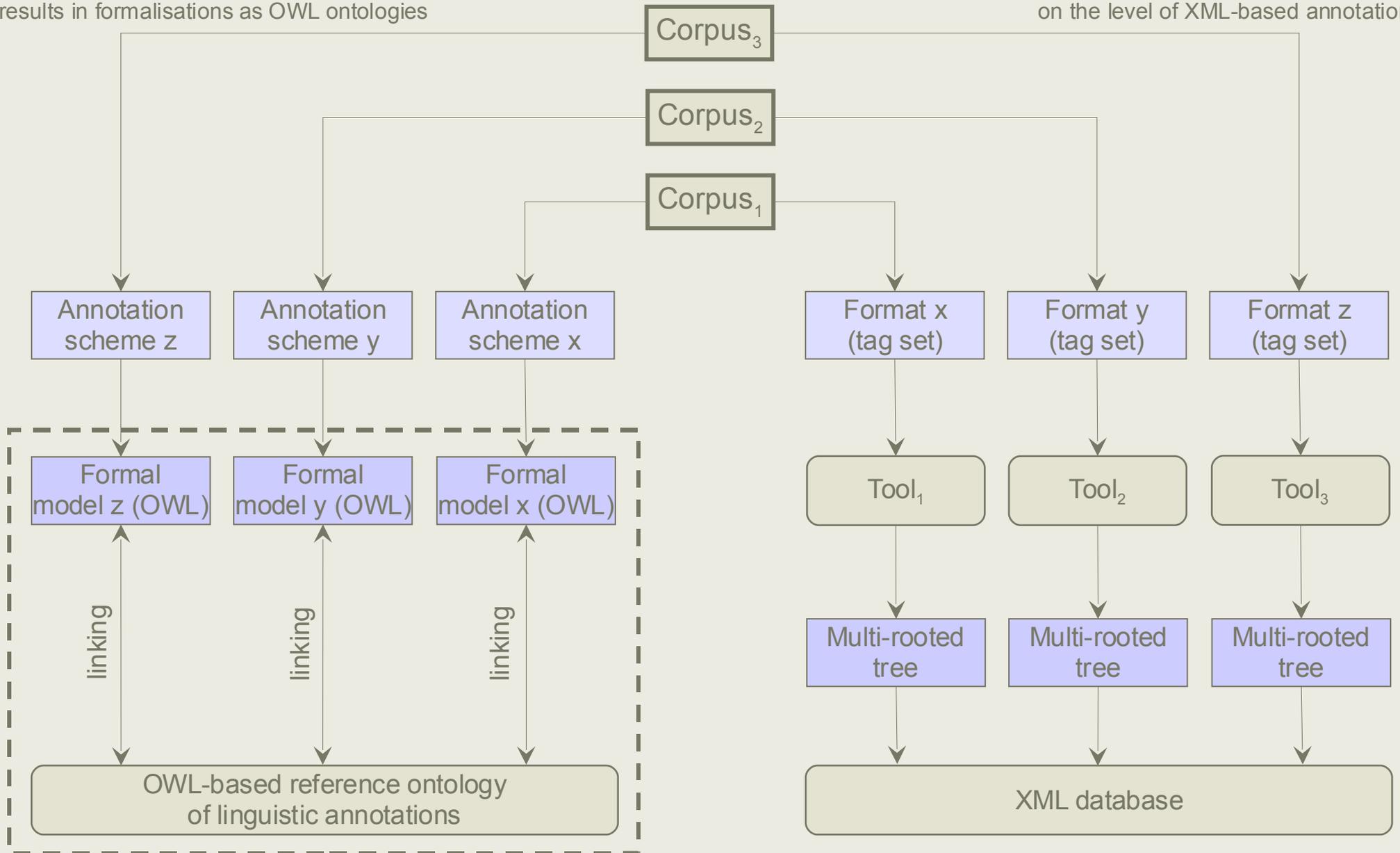
# Processing and Normalisation of Corpus Data

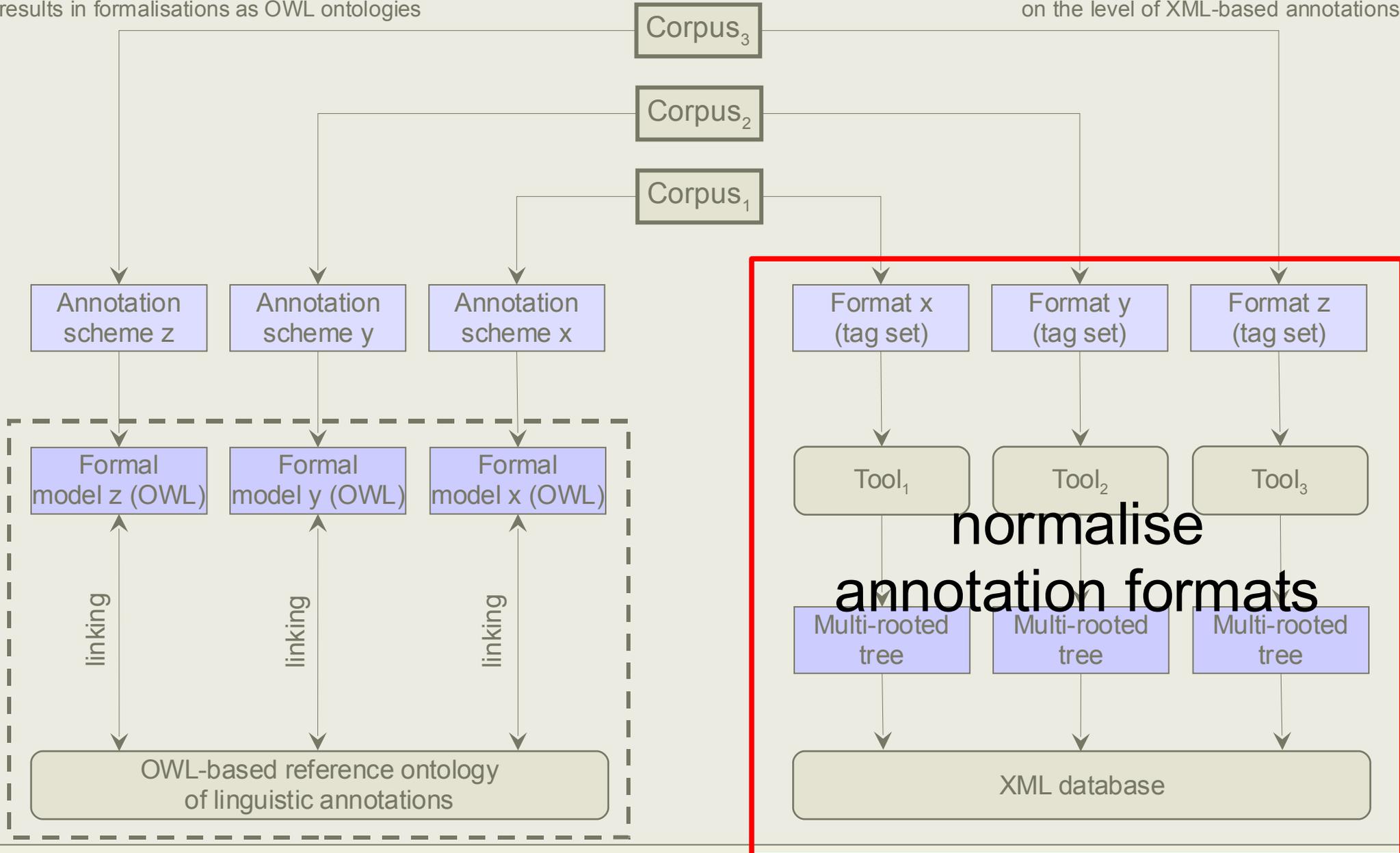Manual analysis of annotation schemes and annotation layers results in formalisations as OWL ontologies

Semi-automatic processing and normalisation on the level of XML-based annotations

$Corpus_3$

$Corpus_2$

$Corpus_1$

| Annotation scheme z | Annotation scheme y | Annotation scheme x | Format x (tag set) | Format y (tag set) | Format z (tag set) |

| Formal model z (OWL) | Formal model y (OWL) | Formal model x (OWL) | $Tool_1$ | $Tool_2$ | $Tool_3$ |

linking    linking    linking

| Multi-rooted tree | Multi-rooted tree | Multi-rooted tree |

OWL-based reference ontology of linguistic annotations

XML database

Ontology-Based XQuery'ing of XML-Encoded Language Resources on Multiple Annotation Layers

# Processing and Normalisation of Corpus Data

Manual analysis of annotation schemes and annotation layers results in formalisations as OWL ontologies

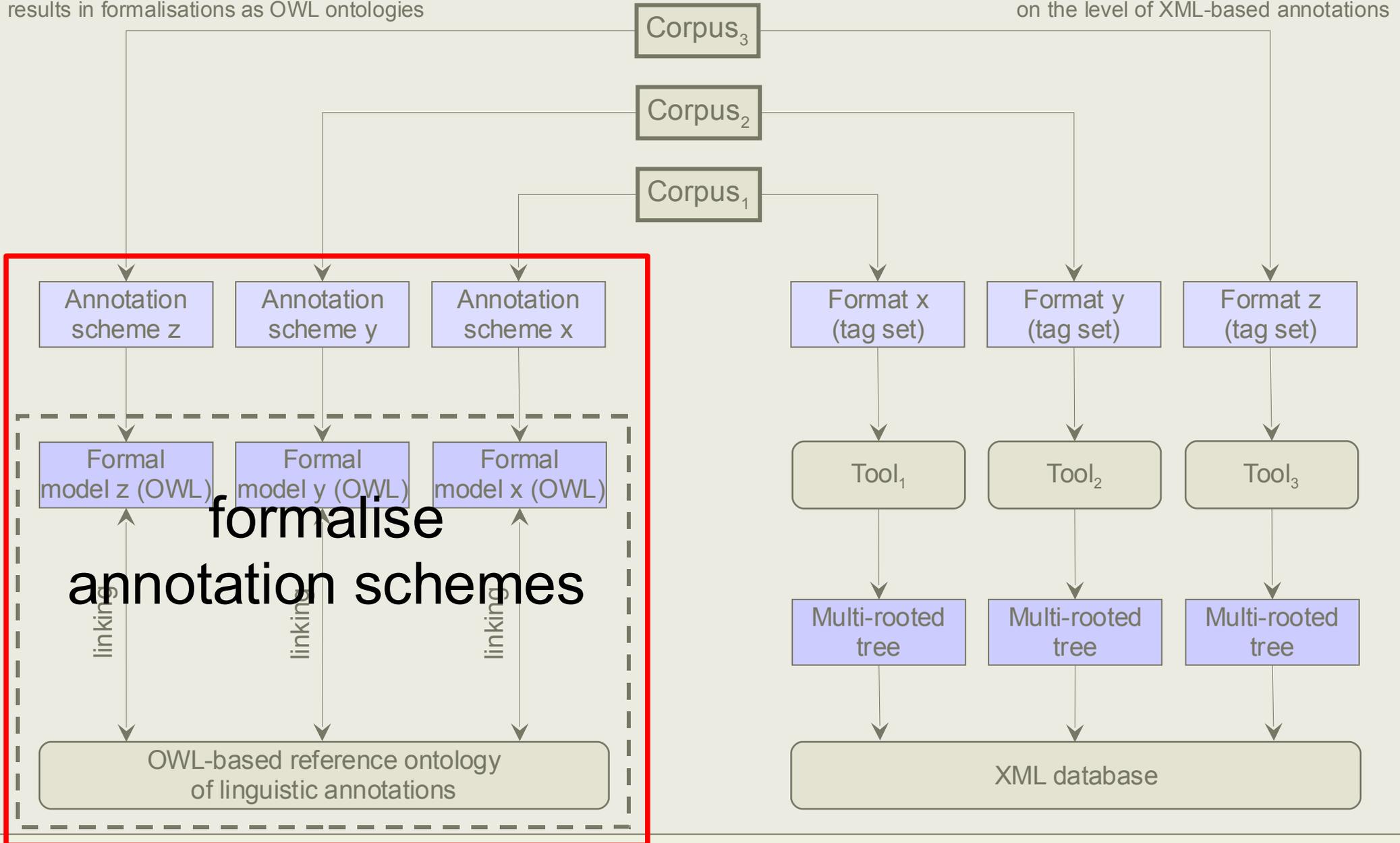Semi-automatic processing and normalisation on the level of XML-based annotations

$Corpus_3$

$Corpus_2$

$Corpus_1$

| Annotation scheme z | Annotation scheme y | Annotation scheme x | | Format x (tag set) | Format y (tag set) | Format z (tag set) |

| Formal model z (OWL) | Formal model y (OWL) | Formal model x (OWL) | | $Tool_1$ | $Tool_2$ | $Tool_3$ |

linking    linking    linking

## normalise annotation formats

| Multi-rooted tree | Multi-rooted tree | Multi-rooted tree |

OWL-based reference ontology of linguistic annotations

XML database

Ontology-Based XQuery'ing of XML-Encoded Language Resources on Multiple Annotation Layers

# Normalising Annotation Format

- Model: multi-rooted trees

- XML-encoded corpora split into multiple layers (trees)

  - One XML file per annotation layer

  - All are identical with regard to their primary data

- Normalizing the XML elements and attributes

  - Tool supported and flexibly configurable (Splitter, Leveler)

- Single layer can be queried with standard XML methods

- Multiple layers cannot be queried with standard methods

  - Introduce custom XQuery functions

# Processing and Normalisation of Corpus Data

Manual analysis of annotation schemes and annotation layers results in formalisations as OWL ontologies

Semi-automatic processing and normalisation on the level of XML-based annotations
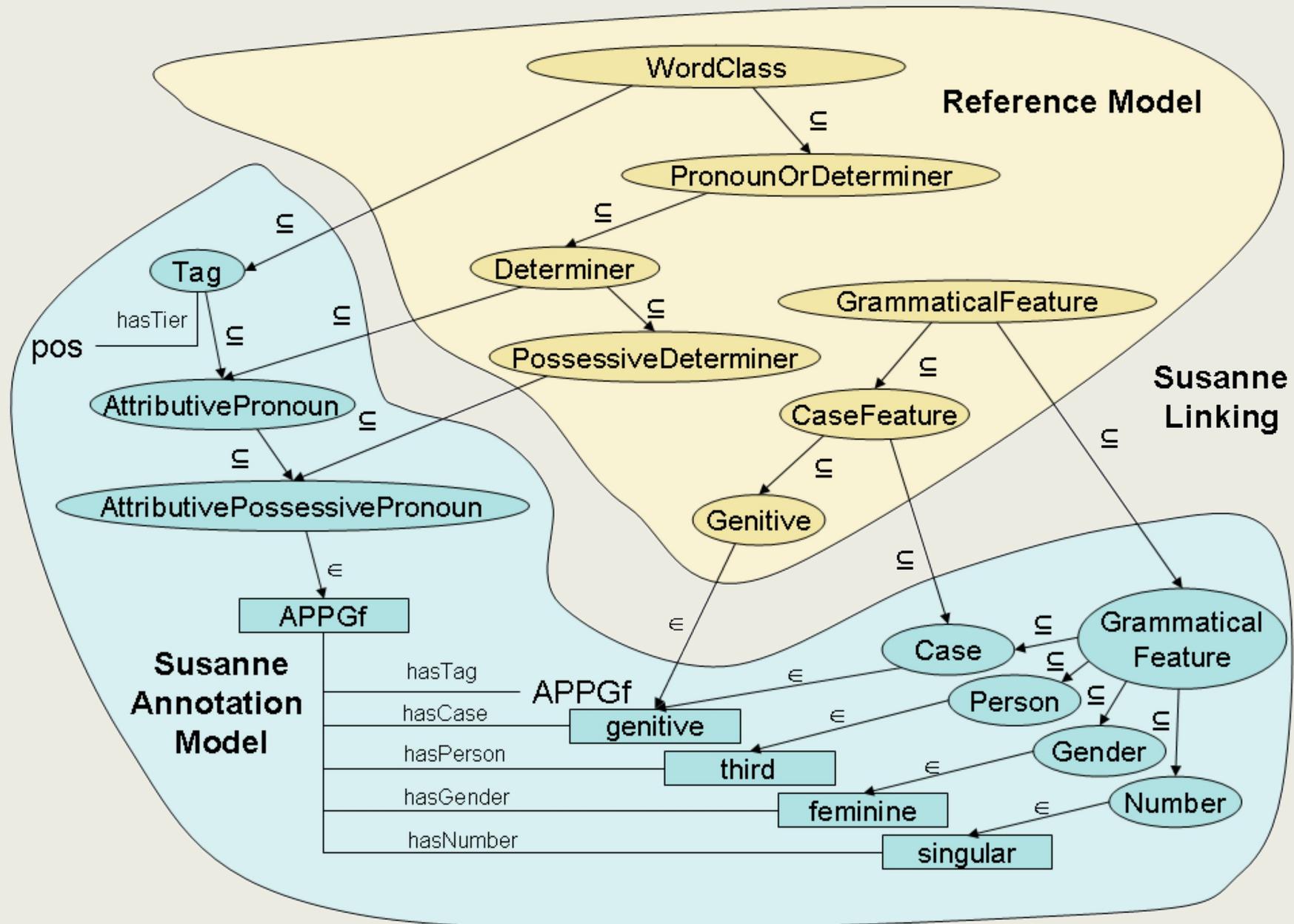
$Corpus_3$

$Corpus_2$

$Corpus_1$

| Annotation scheme z | Annotation scheme y | Annotation scheme x |
|---|---|---|

| Format x (tag set) | Format y (tag set) | Format z (tag set) |
|---|---|---|

| Formal model z (OWL) | Formal model y (OWL) | Formal model x (OWL) |
|---|---|---|

## formalise annotation schemes

linking   linking   linking

$Tool_1$   $Tool_2$   $Tool_3$

| Multi-rooted tree | Multi-rooted tree | Multi-rooted tree |
|---|---|---|

OWL-based reference ontology of linguistic annotations

XML database

Ontology-Based XQuery'ing of XML-Encoded Language Resources on Multiple Annotation Layers

# Formalising Annotation Semantics

- Corpora differ in their annotation schemes

- Integrated treatment of heterogeneous resources requires

  - Annotation specifics documented using a formal language

  - Integrated access to resources with different annotations

- Ontology-based approach

  - Ontological formalisation of annotation schemes

  - Standard format (OWL/DL)

  - Supported by several tools (Protégé, Pellet)
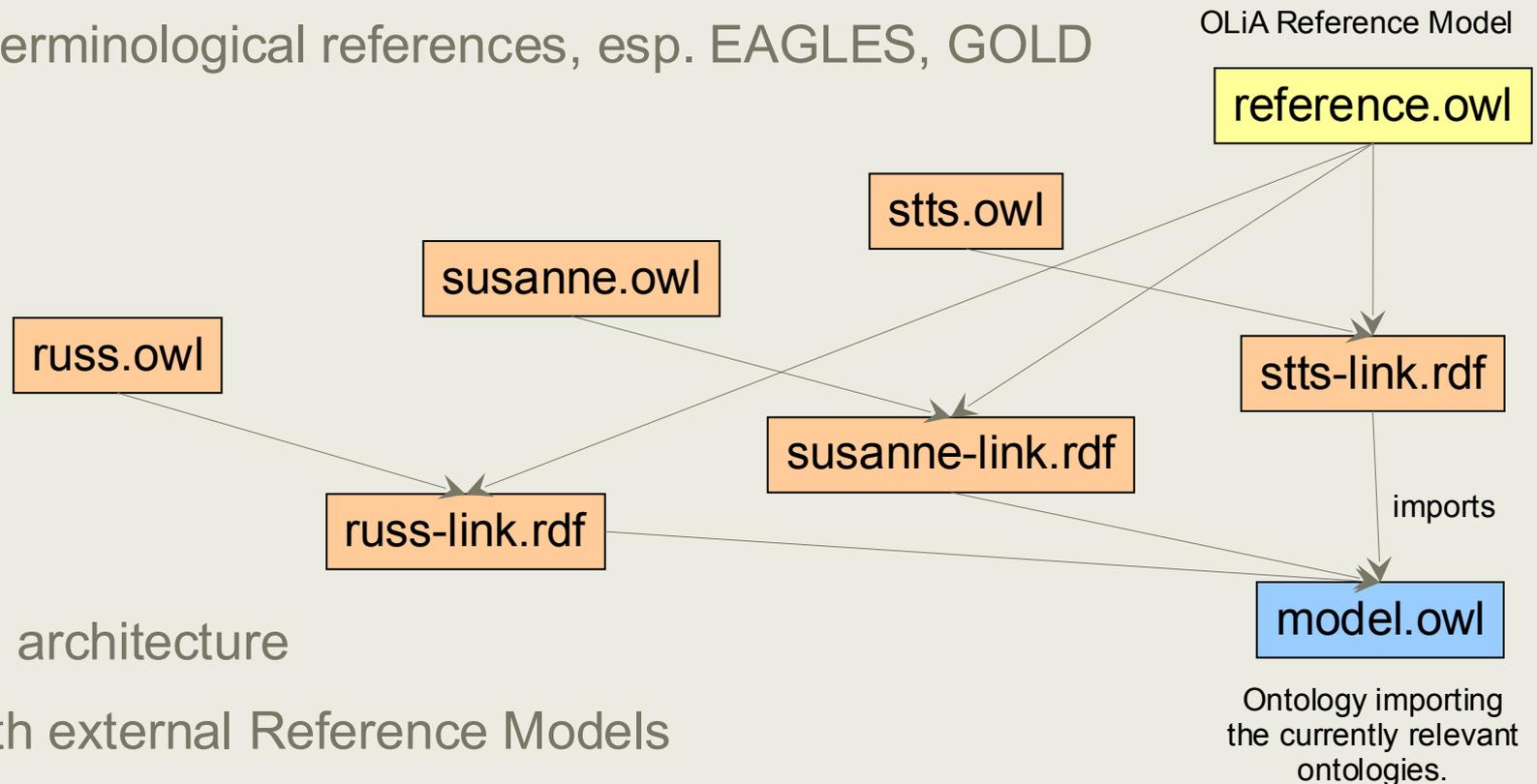
# OLiA: Ontology of Linguistic Annotations

- Annotation Model

  - Ontological formalization of one particular annotation scheme

- OLiA Reference Model

  - Ontological formalization of reference terminology

- Linking

  - Concepts (and tags) of an annotation model are defined with reference to the OLiA Reference Model

    - Sub-concepts/sub-properties     $\subseteq \in$
    - Complex expressions        $\cap \cup$

- An example

  - POS tag APPGf "her" [Susanne Tagset]

# OLiA: Ontology of Linguistic Annotations

# OLiA: Ontology of Linguistic Annotations

- Annotation model
    - 10 models for European and non-European languages
    - POS, morphology, syntactic labels, co-reference, information structure
- OLiA Reference Model
    - Based on terminological references, esp. EAGLES, GOLD

OLiA Reference Model

**reference.owl**

**stts.owl**

**susanne.owl**

**russ.owl**

**stts-link.rdf**

**susanne-link.rdf**

**russ-link.rdf**

imports

**model.owl**

Ontology importing the currently relevant ontologies.

- Linking
    - Extensible architecture
    - Linking with external Reference Models
    - (GOLD, OntoTag, Data Category Registry) supported

# Graphical Query Interface Requirements

- Intuitively usable graphical query interface

- Work with multi-rooted trees

- Include the ontology of linguistic annotations into queries

- Work with open standards, i.e., XQuery, OWL

# SPLICR Graphical Query Interface

- SPLICR has an intuitive graphical query interface

- Generalises over the underlying data structures and querying

- Tree fragment query editor

    - Ontology-supported abstraction of linguistic concepts

    - Operands glue together concepts to construct complex queries

- Multiple display and visualisation modes

    - plain text view                    XML view

    - graphical tree view            time-line view

- Ajax (Asynchronous JavaScript and XML)

- Query and visualisation extensible through modules

# Querying

XML-$1_1$  XML-$1_2$  XML-$1_3$  XML-$1_n$

XML-$2_1$  XML-$2_2$  XML-$2_3$  XML-$2_n$

XML-$n_1$  XML-$n_2$  XML-$n_3$  XML-$n_m$

XQuery engine

XML Database

Ontology

Input (XQuery)

Output (XML)

System database

Intermediate representation

Visualisation

Visualisation

Visualisation

Graphical Query Interface

Free XQuery input

Ontology-Based XQuery'ing of XML-Encoded Language Resources on Multiple Annotation Layers

# Tree Fragment Query Editor



**Corpus Platform Web Access**

Select corpus: [TuBa-D/Z newspaper corpus ▼] [Browse] Search: [Tree Description Search ▼] [Search]

Display results: ○ as list on single page ● with one result per page (verbose)

[□] [□] [↓] [→] [↗]    cat: [VC]    pos: [VAFIN]    tok: [ ]

cat: [VC]
cat: [VCE] [v]

cat: [NX]

any [▼]

HD [▼]

pos: [VAINF]
pos: [VMINF] [v]
pos: [VVINF] [v]

ADJA -- attributive adjective
ADJD -- adverbial or predicative adjective
ADV -- adverb
APPR -- preposition; left circumposition
APPRART -- preposition + article
APPO -- postposition
APZR -- right circumposition
ART -- definite or indefinite article
CARD -- cardinal number
FM -- foreign language material
ITJ -- interjection
KOUI -- subordinating conjunction with "zu" + infinitive

Press 'D' to delete this node.

Ready.

Ontology-Based XQuery'ing of XML-Encoded Language Resources on Multiple Annotation Layers

# Graphical Tree Visualisation



**Corpus Platform Web Access**

Select corpus: [TuBa-D/Z newspaper corpus ▾] [Browse]   Search: [Plain Text Search ▾] [_____] [Search]
Display results: ○ as list on single page ● with one result per page (verbose)

[Previous] Result 5 of 100 [Next]

[TEXT][XML][TREE] [direct ▾] [normal ▾] [Zoom: 75 % ▾]

# AnnoLab Multi-layer Query Example

- Lexical layer  - find the verb *will* ('V')

- Field layer     - find Vorfelds ('VF')

- Coordination - keep those Vorfelds containing *will* as a verb
                    (seq:containing)

```
let $verb := ds:layer('Lexical')//tok
    [starts-with(pos/text,'V')]
    [.//orth = 'will']

let $vf := ds:layer('Field')//ntNode
    [category='VF']

return seq:containing($vf, $verb)
```

TUEBA1: Find the verb *will* in the Vorfeld

# AnnoLab Multi-layer Query Example

- Lexical layer  - find the verb *will* ('V')

- Field layer    - find Vorfelds ('VF')

- Coordination - keep those Vorfelds containing *will* as a verb
  (seq:containing)

```
let $verb := ds:layer('Lexical')//tok
    [starts-with(pos/text,'V')]
    [.//orth = 'will']

let $vf := ds:layer('Field')//ntNode
    [category='VF']

return seq:containing($vf, $verb)
```

TUEBA2: Find the verb *will* in the Vorfeld

# AnnoLab Multi-layer Query Example using OLiA

- Lexical layer - find the verb *will* ('V')

- Field layer - find Vorfelds ('VF')

- Coordination - keep those Vorfelds containing *will* as a verb
  (seq:containing)

```
let $verb := ds:layer('Lexical')//tok
    [pos/text = oc:expand('Verb')]
    [.//orth = 'will']

let $vf := ds:layer('Field')//ntNode
    [category='VF']

return seq:containing($vf, $verb)
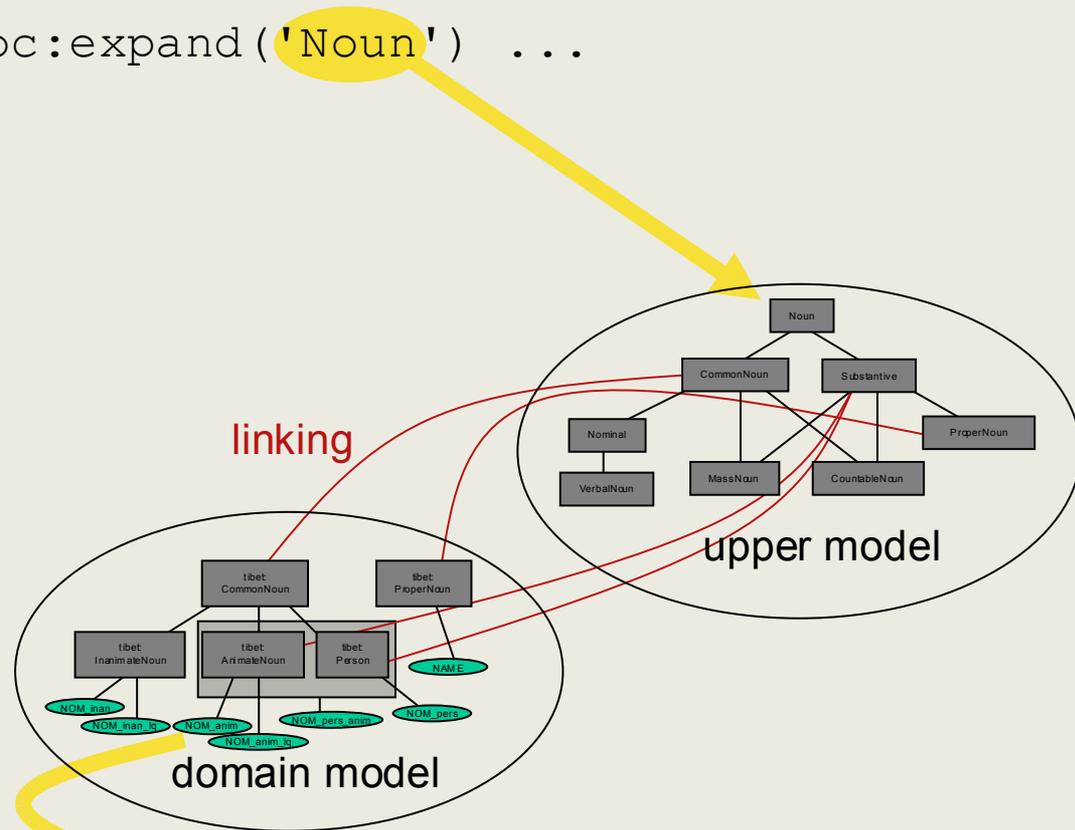```

TUEBA2: Find the verb *will* in the Vorfeld using OLiA

# oc:expand in Detail

corpus query     `... oc:expand(`'Noun'`) ...`

ontology lookup:
1. instance retrieval
2. application of set
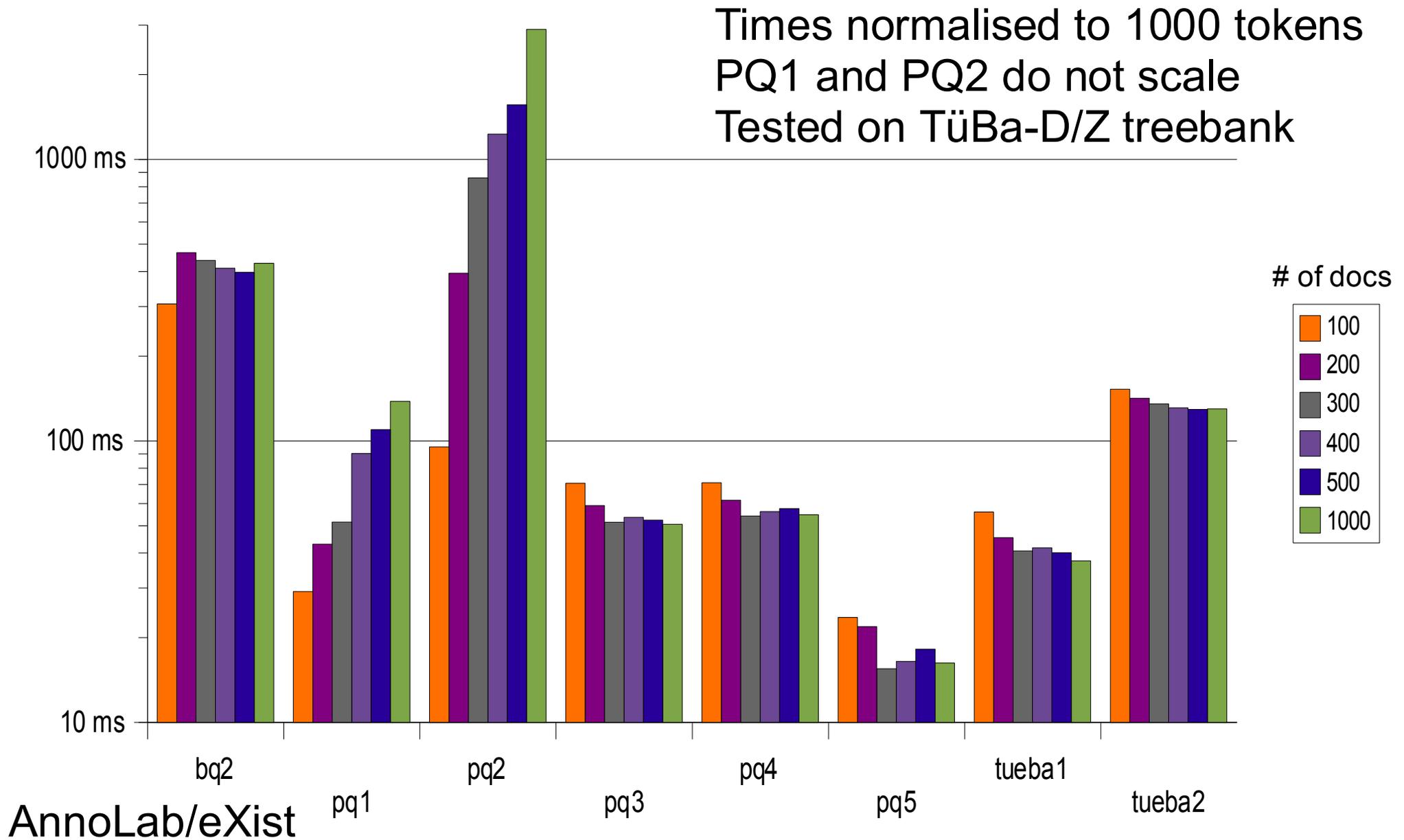    operators



linking

upper model

domain model

`... NOM_inan | NOM_inan_lq | NOM_anim | NOM_anim_lq | NOM_anim_pers | NOM_pers | NAME ...`
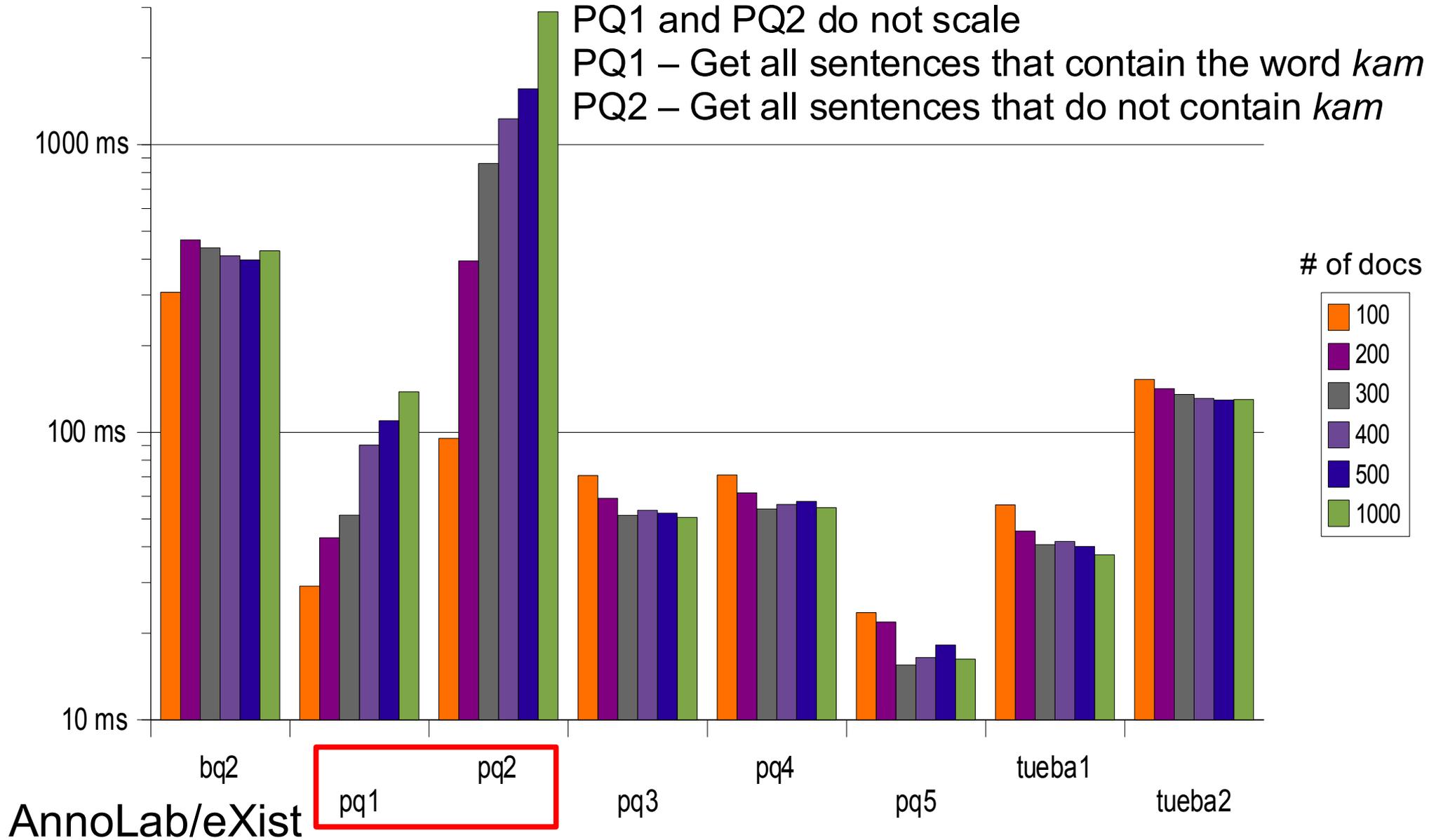
# Experimentation queries

- PQ1 – Get all sentences that contain the word *kam*

- PQ2 – Get all sentences that do not contain *kam*

- PQ3 – Get references to all NPs

- PQ4 – Get all subtrees dominated by NPs

- PQ5 – Get all NPs subtrees dominated by a VP

- TUEBA1 – Find all occurrences of the verb *will* in the Vorfeld

- TUEBA2 – TUEBA1 using OLiA

- BQ2 – Get NPs that are immediate following siblings of a verb

# Average Query Run-Time (logarithmic)



Times normalised to 1000 tokens
PQ1 and PQ2 do not scale
Tested on TüBa-D/Z treebank

# of docs
- 100
- 200
- 300
- 400
- 500
- 1000

1000 ms

100 ms

10 ms

bq2 · pq1 · pq2 · pq3 · pq4 · pq5 · tueba1 · tueba2

AnnoLab/eXist

# Average Query Run-Time (logarithmic)



PQ1 and PQ2 do not scale
PQ1 – Get all sentences that contain the word *kam*
PQ2 – Get all sentences that do not contain *kam*

# of docs

- 100
- 200
- 300
- 400
- 500
- 1000

1000 ms

100 ms

10 ms

bq2

pq1    pq2

pq3

pq4

pq5

tueba1

tueba2

AnnoLab/eXist

# Summary

- Approach to querying XML-annotated corpora using standard techniques such as XPath and XQuery

- Extended an XML database to query multi-rooted trees

- Built an OWL ontology of linguistic annotations generalising over annotation schemes and tag sets

- OWL ontology can be used for query expansion

- Implemented an intuitive and flexible graphical query interface

# Conclusions and Future Work

- Work on SPLICR is ongoing

- Building the GUI to explore and to query meta-data

- Extended query interface functionality (e.g. saved searches)

- Working on benchmark queries for evaluating XML databases with respect to linguistic corpora