

**dboe@TEI**

# **Remodelling a Database of Dialects into a Rich LOD Resource**

**Jack Bowers, Daniel Schopper, Eveline Wandl-Vogt**

Presented at TEI Members Meeting  
October 2015  
Lyon, France

# Database of Bavarian dialects of Austria (DBOE)

- Project collection of vocabulary from Bavarian dialects of what was then the Hapsburg monarchy/Austro-Hungarian Empire began in 1911...
- Estimated 4 million paper slips containing lexical, geographical and other information ....
- Collection of dialectal language data from locations around the former Austro-Hungarian Empire

1911- **Collection begins**



1963 - dictionary (**WBÖ**)

1998 database (**DBÖ**)

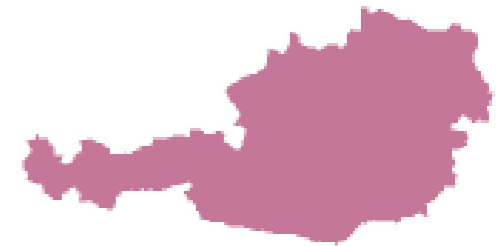
2010 [dbo@ema](mailto:dbo@ema) [online](#)

2012 [wboe<sup>d</sup> online](#)

1998

2013 **Collection ends**

[wboe@SKOS](#)  
(OpenUp! NHM)



# Overview of Scenario:

## Task:

- Convert our legacy dialectal, language data and metadata into formats that are:
  - compatible with one another (SQL and TUSTEP data);
  - ....
  - Make LOD compatible;
  - Bring formatting in line with standards for best practice

## Data:

- **TUSTEP Database:** 2 299 608 single data entry files; (728 drawers)
- MySQL (dbo@ema): +- 200,000 entries (4 drawers)

- Phonetic representation(s) of dialectal forms;
- Grammatical info:pos, number, case, etc.
- Details of word formation;
- Etymological information;
- Translation (&/or) definition of meaning;
- Related forms in other dialectal varieties;
- Location data for each record;
- Questionnaires used in elicitation of dialectal data;
- Bibliographic references;

## Motivations:

- LOD offers means to base semantics for a given concept in predefined, inter/multi-lingual knowledge hubs;
- Desire to have access to dialectal data according to as many different lexical, semantic, geographic, temporal, features as possible;

## Sources & Services:

- dbpedia; wikidata; geoNames;
- NHM (OpenUp!) Data (plantnames)
- BableNet/Babelfy;

## Applicable to:

- Semantics/sense info for language/dialectal data;
- Geographical data;
- Grammatical, lexicographic data (ISOcat)
- Fragebogen; (questionnaires)

# Major Challenges in Converting DBOE to TEI

- **Heterogeneous formatting in TUSTEP records;**

*Dates can appear in a number of formats, with different levels of specificity:*

*e.g. (1936), c(1906), (1913-23)...*

*...many of which are incomplete:*

*e.g. (19xx), (19xx), c(1943-19xx)*

**Some entry records:**

*have empty fields;*

*have same category of data in different places*

\*\*\*\*\*

\*A\* HK 507, k5070426.foe, **unkorr.**

\*HL\* (r<ot)kopfecht:2

\*NL\* :

\*VL\* :

\*QU\* Tiersee Juffinger

\*QN\* {1C.3b03} nUInnggeb.:nwNOTir.:öNTir.

\*^@ FbB.JUFFINGER. (19xx)

[SFb./EFb./Mtlg.] \*O\* Thiers. NTir.

===

\*NR\* 53C5: rot in Komp. (rotfuchset,

Rotblättlein, feuerrot)

\*LT1\* rEotkopfa [ ]

\*BD/LT1\* -- \*ANMO\* von gutaussehenden

Menschen

===

\*\*\*\*\*

\*A\* HK 191, t191^#249.52 =

t1910201.pla^#191.660

\*HL\* tschga:9

\*QU\* Schüttelkopf, Dt. TierN. in Kä., In:

Carinthia (1906) \*S\* II, 66 \*N\* 96

\*QN\* {wb} mkPubl./MdaWs.:TierN \*^@

SCHÜTTELKOPF. (1906) S. [S-5637] \*O\*

===

\*KT1\* Tschgå, tschgå \*KL\* Fg:99 \*O\* Dtl.

\*BD/KT1\* Lockruf für Kälber

\*\*\*\*\*

\*A\* HK 204, d204^#1779.1 = tut1028.res^#1.398

\*HL\* Tschuffit:1

\*QU\* STir. PflN und TierN (Sammelbel.)

^@^@^@

===

\*LT1\* Tschuffit \*ANM\* Duregger \*O\* Meran

\*LT2\* Tschuffitt \*O\* Kardaun

\*BD/LT1\* Zwergohreule; Ephiattes scops Prag

# Major Challenges in Converting DBOE to TEI

- **Invalid (and invisible) non-unicode characters in text out put of both TUSTEP and SQL;**

Plâsepl:□me

F:□tterpl:□me



# Major Challenges in Converting DBOE to TEI

- **Labeling @xml:lang for dialectal data;**

BCP 47 only specifies to the degree of country & administrative region;  
*(language)-(country)-(region);*  
*(ISO639)-(ISO3166)-(3166-2)*

*There is only perhaps one case when this is specific enough;*

***xml:lang="bav-AT-9"      Vienna (ok)***

*However,*

***xml:lang="bav-AT-7"      Tyrol (too vague)***

*Speech of places in the same region very often are very different, thus we need more specifics:*

***xml:lang="bav-AT-7-zillertal"      (not valid)***

*But... individual towns, place names, or any other geographical identifiers do not have official status in language identification systems*

***xml:lang="bav-AT-7-x-zillertal"      (private)***



# Major Challenges in Converting DBOE to TEI



ÖSTERREICHISCHE  
AKADEMIE DER  
WISSENSCHAFTEN

- **Phonetic transcriptions not in reliably consistent notation:**
  - *Individual data collectors may or may not have followed instructions to use common phonetic notation (Wienerisch Teuthonista)*
  - *Different conventions were often used by data scientists entering the handwritten transcriptions into the DB;*
  - *Representation of phonetic transcriptions in DB based on the notation of the handwritten transcription and not the phonetic forms themselves!*

*Thus, the accuracy of any attempt to convert data to IPA cannot be guaranteed!*

# Major Challenges in Converting DBOE to TEI

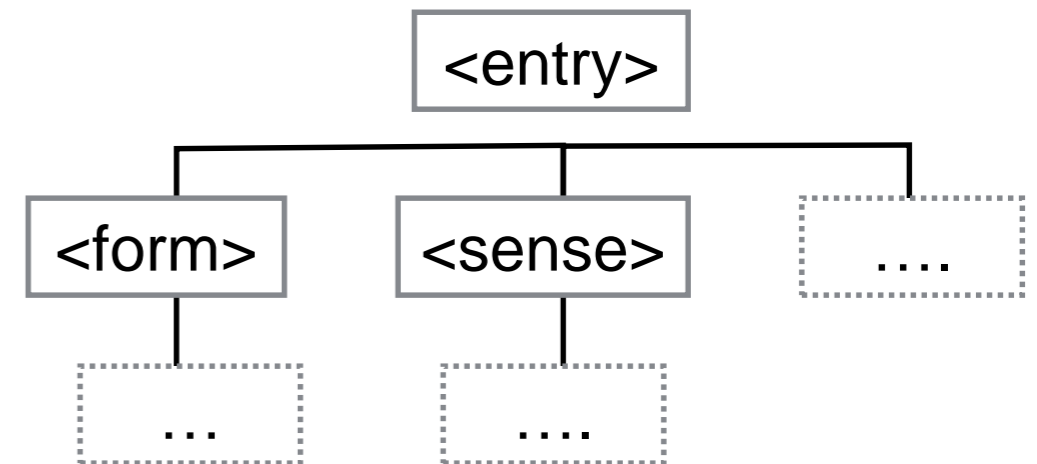
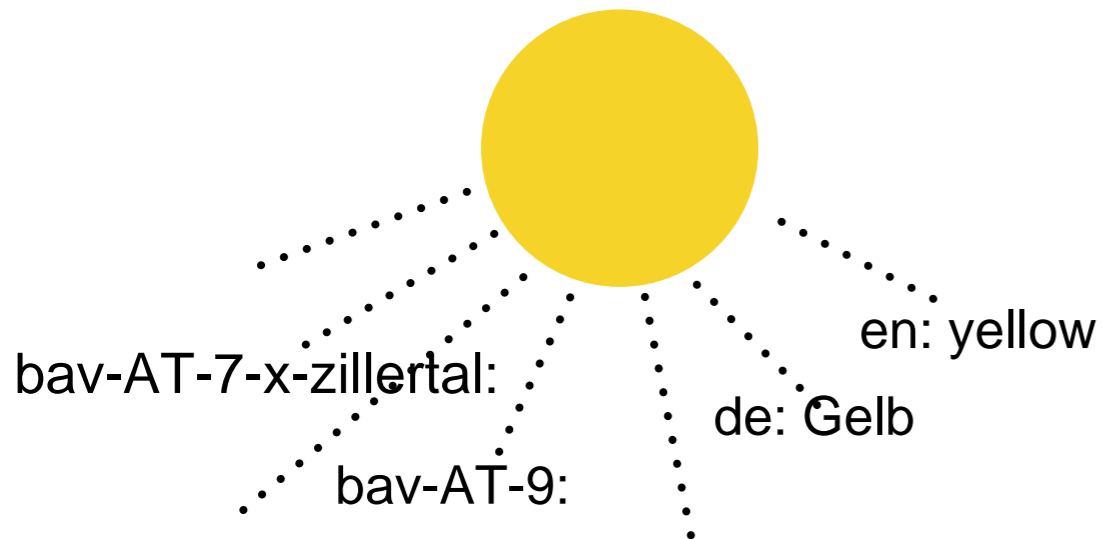
- **What is the best type of TEI output for data?**

## *Onomasiological or Semisiological*

Concept -> term

Form -> Sense

TEI Dictionary?



- **should we consider other formats?**  
(TBX-TEI, hybrid?)

# Potential Output: TEI Dictionary with Etymology

dialect-form: **dotEjrgöjlb**

de: **dottergelb**

etymological-form :**(dotter)gëlb:2**

dialect-place :**Gm. Wien**

```
<entry xml:id="dottergelb-adj-wien-gersthof">
  <form type="lemma">
    <pron notation="tustep" xml:lang="ba-AT-9">
      <seg xml:id="dtrGlbwnGer01">dotEj</seg>
      <seg xml:id="dtrGlbwnGer02">röjlb</seg>
    </pron>
    <gramGrp>
      <pos>adj</pos>
    </gramGrp>
  </form>
  <sense>
    <usg type="dom" corresp="http://dbpedia.org/resource/Color">Color</usg>
    <usg type="geo">
      <placeName type="district" ref="place:W">Wien-Gersthof</placeName>
    </usg>
    <cit type="translation">
      <oRef xml:lang="de">dottergelb</oRef>
    </cit>
    <etym type="compounding">
      <etym type="metonymy">
        <cit type="etymon" corresp="#dotter-Subs-wien">
          <pRef corresp="#dtrGlbwnGer01">dotEj</pRef>
        </cit>
      </etym>
      <cit type="etymon" corresp="#gelb-adj-wien">
        <pRef corresp="#dtrGlbwnGer02">röjlb</pRef>
      </cit>
    </etym>
  </sense>
</entry>
```

*Formatting in accordance to Bowers & Romary (Upcoming)*

# Workflow: TUSTEP > XML > TEI

## TUSTEP entry

\*A\* HK 288, F288^#115.1 = F2880826.eck^#72.1  
\*HL\* Fuchs:1  
\*QU\* Schlagen b. Gmunden, Loitlesberger  
\*QN\* {5.3a06} nöSkg.:swTraunv.:OÖ \*^@ Slg.LOITLESBERGER.  
(19xx) [Slg."PflN+TierN.Gmunden"+Erg.] \*O\* Schlagen in Gm.  
Gmunden OÖ  
===  
\*NR\* 53C6: rot versch. Abstufg. (purpurn, der Purpur)  
\*LT1\* Fuks  
\*BD/LT1\* rotes oder rötliches Pferd

TEI <fs> (grammatical features) TEI (listPlace)



## TEI(list)

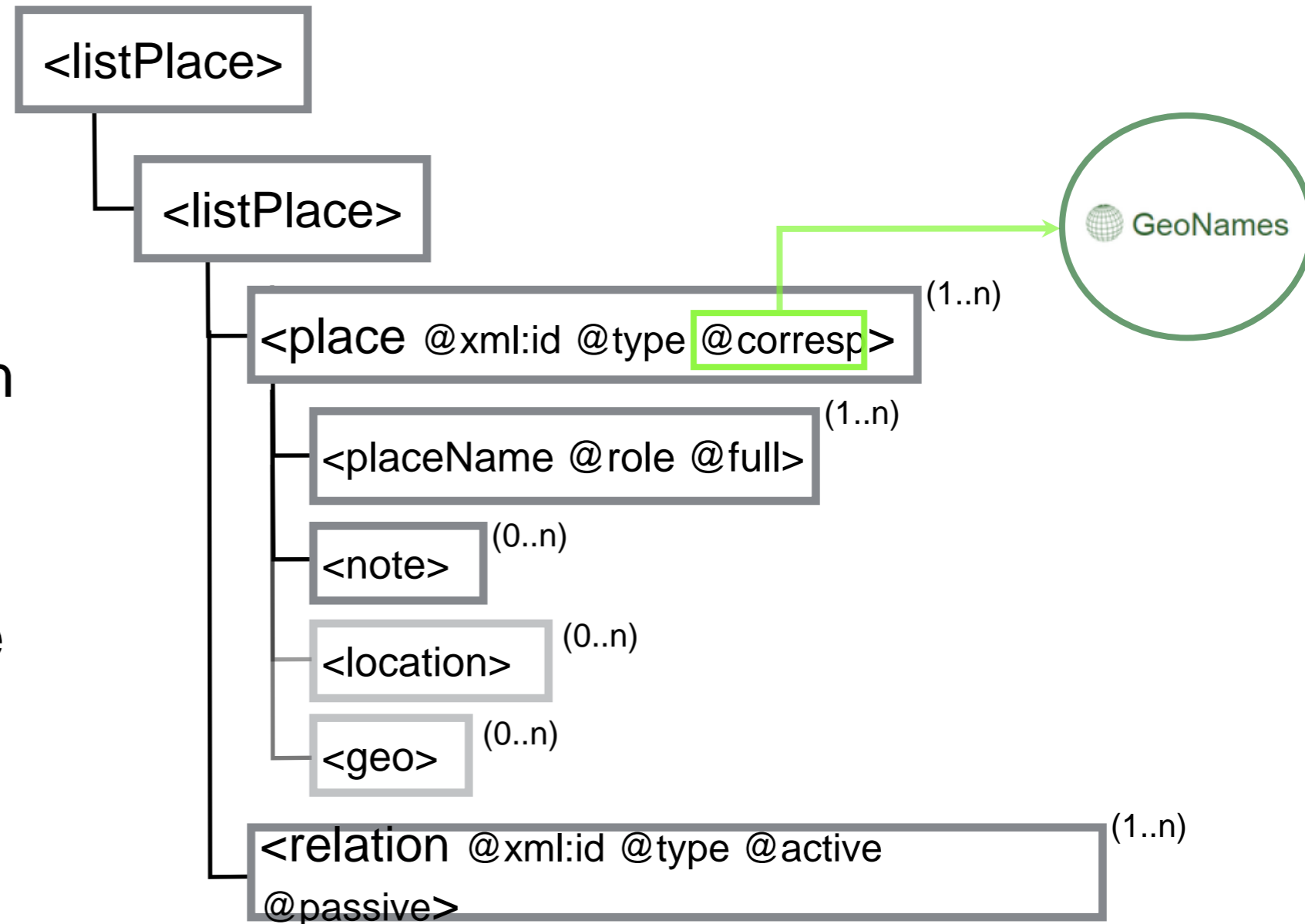
```
<item xml:id="d1e437" n="21">  
  <ref type="archive">HK 288, F288^#115.1 = F2880826.eck^#72.1</ref>  
  <cit type="etymon">  
    <pRef>Fuchs</pRef>  
    <gramGrp>  
      <pos sameAs="lemma_wortat-tei-fs.xml#Sub">Sub</pos>  
    </gramGrp>  
  </cit>  
  <ref type="source">Schlagen b. Gmunden, Loitlesberger</ref>  
  <bibl>  
    <ref>5.3a06</ref> nöSkg.:swTraunv.:OÖ  
    Slg.LOITLESBERGER. (<date>19xx</date>  
    [Slg."PflN+TierN.Gmunden"+Erg.]</bibl>  
  <placeName>Schlagen in Gm. Gmunden OÖ</placeName>  
  <ref type="questionnaire">53C6: rot versch. Abstufg. (purpurn, der Purpur)</ref>  
  <cit type="lautung" xml:id="d1e451" n="1">  
    <pRef>Fuks</pRef>  
  </cit>  
  <cit type="translation" corresp="#d1e451">  
    <oRef xml:lang="de">rotes oder rötliches Pferd</oRef>  
  </cit>  
</item>
```

## BASIC XML

```
<record n="21">  
  <A>HK 288, F288^#115.1 = F2880826.eck^#72.1</A>  
  <HL>Fuchs:1</HL>  
  <QU>Schlagen b. Gmunden, Loitlesberger</QU>  
  <QN>  
    <bibl>{5.3a06} nöSkg.:swTraunv.:OÖ *^@  
    Slg.LOITLESBERGER. (19xx)  
    [Slg."PflN+TierN.Gmunden"+Erg.]</bibl>  
    <O>Schlagen in Gm. Gmunden OÖ</O>  
  </QN>  
  <NR>53C6: rot versch. Abstufg. (purpurn, der Purpur)</NR>  
  <LT1>Fuks</LT1>  
  <BD_LT1>rotes oder rötliches Pferd</BD_LT1>  
  <orig>....</orig>  
</record>
```

# Key TEI Documents: listPlace

- Inventory of all known places from which we have language data;
- Establishing inventory of such references enables the option of tagging with attribute or element values in using data;
- Reusable and extensible w/in our project and institution;
- Enhance original information from LOD sources like GeoNames..



# Defining Location: Relational Hierarchies

Österreich

Tirol

Nord Tirol

Mittel Nord Tirol

Zillertal

```
...
<place xml:id="Öst" type="region" subtype="administrative"
  corresp="http://sws.geonames.org/2782113/about.rdf">
  <placeName role="main">
    <placeName full="abb" xml:lang="de">Öst.</placeName>
    <placeName full="yes" xml:lang="de">Österreich</placeName>
  </placeName>
  <location>
    <geo corresp="#gis_region_id-862"/>
  </location>
</place>
```

```
...
<relation name="parentRegion" active="#Öst" passive="#Bgl #Kä #NÖ #OÖ #Sa #St #Tir #W #OÖst #SÖst #NÖst #WÖst"/> *
```

```
<relation name="parentRegion" active="#Tir" passive="#NTir #OTir #TirHocht #sbairTir #TirÜGeb"/>
```

```
<relation name="parentRegion" active="#NTir" passive="#wNTir #öNTir #mNTir #Innt"/>
```

```
** <relation name="parentRegion" active="#mNTir" passive="#Stub #Zillert #Sillgeb #mInntInnsbr"/>
```

```
*** <place xml:id="Zillert" type="region" corresp="http://sws.geonames.org/2760567">
  <placeName role="main">
    <placeName full="abb" xml:lang="de">Zillert.</placeName>
    <placeName full="yes" xml:lang="de">Zillertal</placeName>
  </placeName>
  <location> (i..n)
    <geo corresp="#gis_region_id-771"/>
  </location>
  <note>häufige Abk. "Zill.": Tal der Ziller samt einmündender Alpentäler im SO vom —: #mNTir.# --
  - vom Sammler Michael Juffinger mit der Ziffer IV bezeichnet ---</note>
</place>
```



# Defining Location: LOD Aspect

de: Zillertal

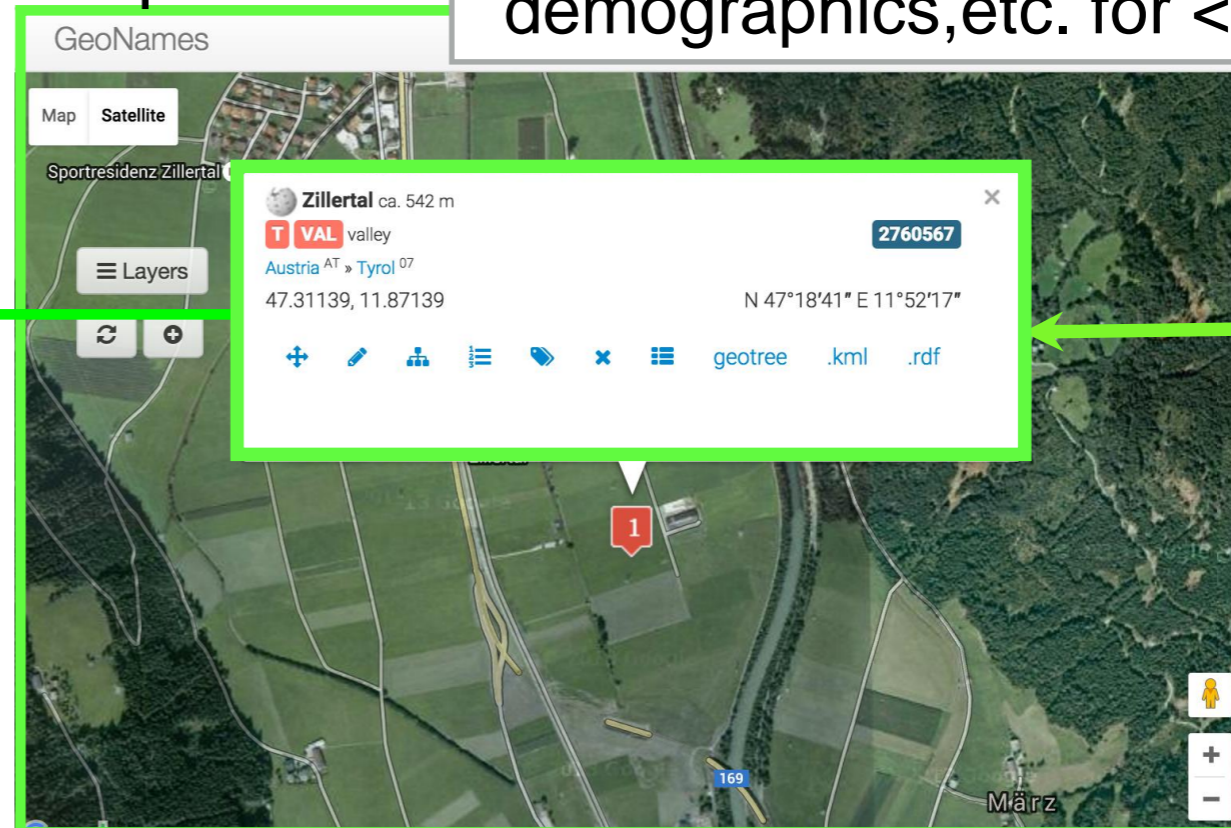
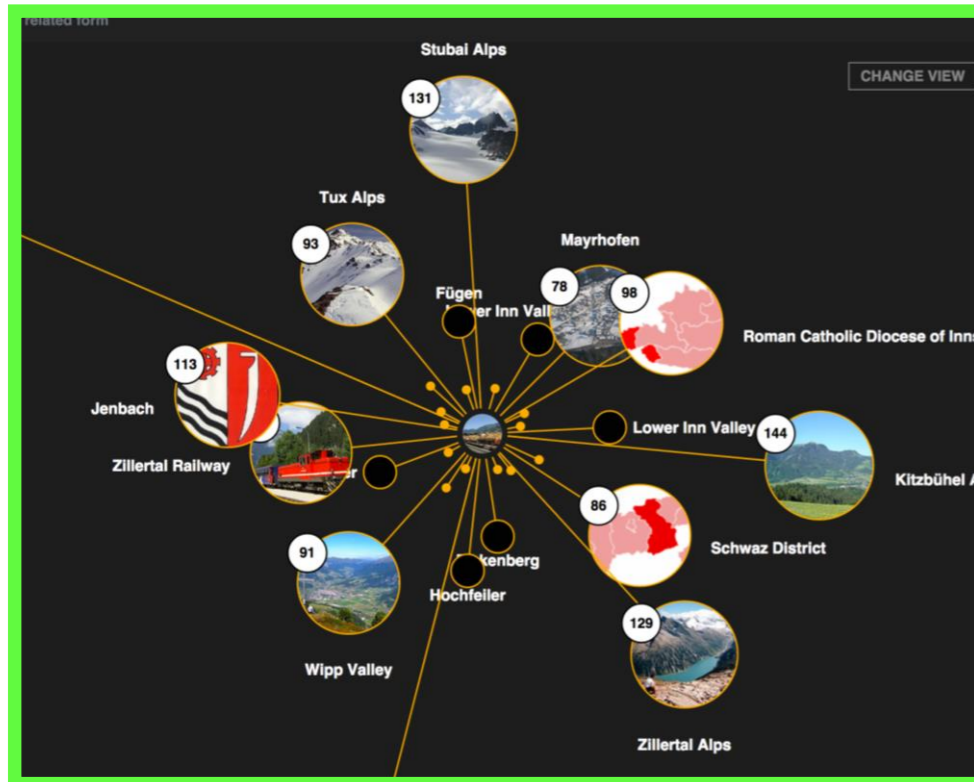
TEI <place>

geoNames uri

```
<place xml:id="Zillert" type="city" corresp="http://sws.geonames.org/2760567">  
<placeName role="main">  
  <placeName full="abb" xml:lang="de">Zillert.</placeName>  
  <placeName full="yes" xml:lang="de">Zillertal</placeName>  
</placeName>  
<location>  
  <geo corresp="#gis_region_id-771"/>  
</location>  
<note>häufige Abk. "Zill.": Tal der Ziller samt einmündender Alpentäler im SO vom ---: #mNTir.#  
--- vom Sammler Michael Juffinger mit der Ziffer IV bezeichnet ---</note>  
</place>
```

geoNames: linked data, related concepts

geoNames: geo co-ordinates,  
demographics, etc. for <place>





# Defining Location: LOD Aspect

Place:

en: Tyrol  
de: Tirol

TEI <place>

```

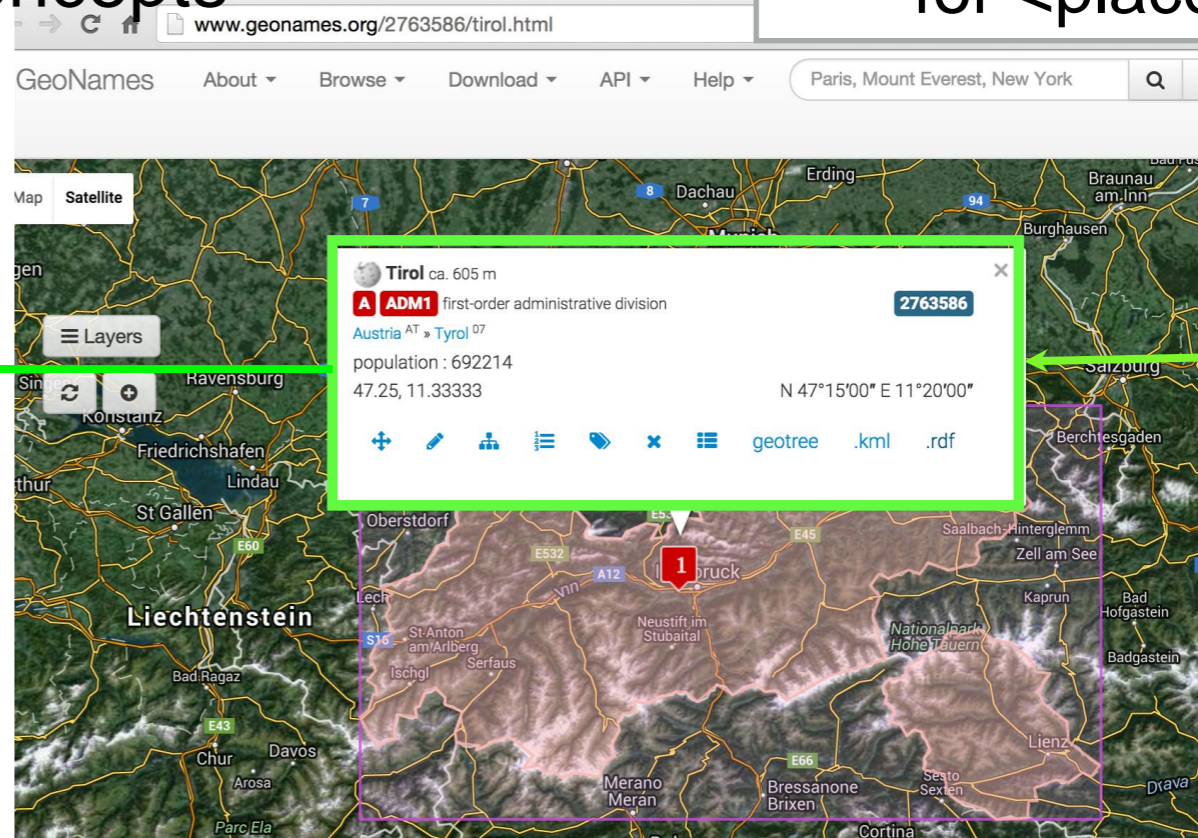
...
<place xml:id="Tir" type="region" subtype="administrative"
  <placeName role="main">
    <placeName full="abb" xml:lang="de">Tir.</placeName>
    <placeName full="yes" xml:lang="de">Tirol</placeName>
  </placeName>
  <location>
    <geo corresp="#gis_region_id-762"/>
  </location>
  <note>nicht eindeutige und daher eher zu vermeidende Bez.: "Alttirol" für --:
    #OTir.#, #NTir.#, #STir.# und Bld. Tirol mit --: #OTir.#, #NTir.# (EKü.)</note>
</place>
    
```

geoNames uri

<http://sws.geonames.org/2764958>

geoNames: geo co-ordinates, demographics, etc. for <place>

geoNames: linked data, related concepts

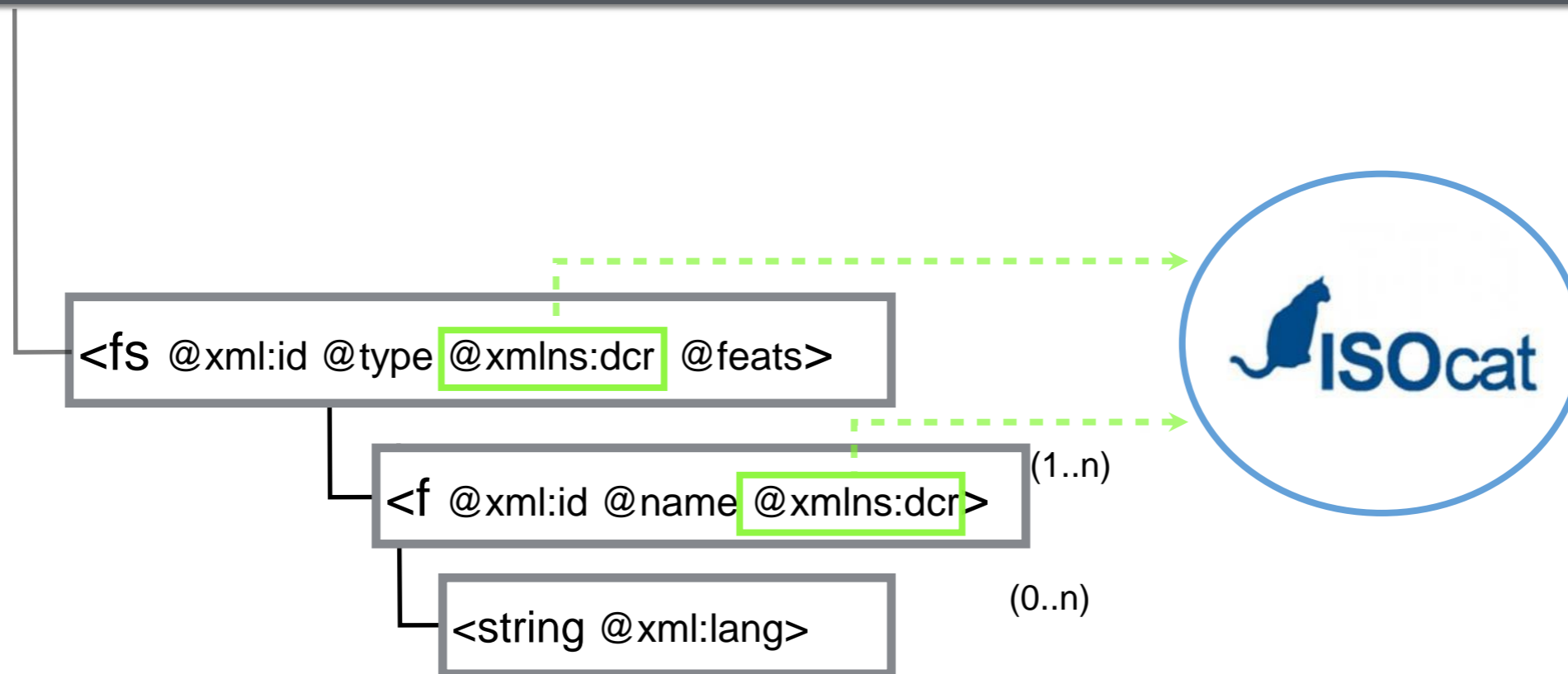


# Key TEI Documents: Feature Structures for Lexical, Grammatical Categories



- Original (TUSTEP) entries has extensive inventory of grammatical and lexical features;
- TEI feature structures compiled manually according to documentation of guidelines;
- Basic part of speech tag (subs, verb, adj, adv, interj, conj,...etc) inherited from etymological descendant form (Hauptlemma);
- Further detail as to subcategorization, context, inflection form(s), and collocation specified for each dialectal form (Lautung);

# Key TEI Documents: Feature Structures for Lexical, Grammatical Categories



`<fs>` used for categorical features, and for complex categories (needing to reference combine multiple features);

`<f>` used for most basic categories and for features for which we want to specify value in one or more strings;

# Grammatical information in dialectal data

e.g.

**\*HL\* :1 = Substative**

```
<f name="Substantiv" n="1" xml:id="Sub"  
xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-3347">  
....  
</f>  
...  
<fs type="Gender" xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-3217">  
....  
<f name="MasculineGender" xml:id="Masc"  
xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-3312">  
  <string>m</string>  
</f>  
....  
</fs>
```

\*A\* HK 165, d165A#1993.1 = T1650708.sch#18

**\*HL\* T<ötling:1**

\*QU\* Schüttelkopf Dt. Tiern in Kä. - In: Carinthia II,96 (1906) \*S\* 66

\*QN\* {wb} mkPubl./MdaWs.:TierN \*^@ SCHÜTTELKOPF. (1906) S.  
[S-5637] \*O\* ob. Mllt.

===

**\*LT1\* Tötling [m]**

\*BD/LT1\* abgESPäntes Kalb

**\*LT1\* [m] = Masculine**



# Grammatical information in dialectal data

e.g.

**\*HL\* :1 = Substative**

```
<f name="Substantiv" n="1" xml:id="Sub" xmlns:dcr="http://www.isocat.org/ns/dcr"
dcr:datcat="http://www.isocat.org/datcat/DC-3347">
....
<fs type="diminutive" xml:id="SubD" dcr:datcat="http://www.isocat.org/datcat/DC-4922">
  <!-- D = Deminutiv/Diminutive = Verkleinerungsform/Koseform (ohne näh.
  Unterscheidung);-->

  <f name="diminutive-lein" xml:id="SubD1"/>
  <!-- D1 = 1. Deminutiv: Häusl, Dearfö (Dörflein), Stüble (lemmatisiert -lein)-->

  <f name="diminutive-ellein" xml:id="SubD2"/>
  <!-- D2 = 2. Deminutiv: Häuserl, Hansei, Hantale, Fußile (lemmatisiert -ellein) -->

  <f name="diminutive-i" xml:id="SubD3"/>
  <!-- D3 = Koseform auf -i: Greti, Pferti, Katzi -->
</fs>
....
```

\*A\* HK 295, f295^#2793.1 = f2951202.pir^#108.1

\*HL\* **Procken:1**

\*QU\* Tullnitz Mühlhauser

\*QN\* {8.3c04} Jaispitzt.:nöSMä. \*^@ FbB.MÜHLHAUSER. (1922-35)  
[SFb.1-6,12-108/EFb.4-9/WrFb.1/Mtlg.] \*O\* Tullnitz SMä.

===

\*NR\* 30C23 (F): gr/kl. Brotstück (Reankn, Bröckel,\*); Ra.; Komp.

\*LT1\* **Brèkl [D1]**

\*LT2\* **Brèkajrl [D2]**

\*BD/LT1\* kleines und kleinstes Stückchen Brot

**\*LT1\* [D1] = Diminutive (dialectal cognate to High German ‘-lein’)**

**\*LT2\* [D2] = Diminutive (dialectal cognate to High German ‘-ellein’)**

# Grammatical information in dialectal data

e.g.

**\*HL\* :1 = Substative**

```
<f name="Substantiv" n="1" xml:id="Sub" xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-3347">  
</f>  
<fs type="Gender" xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-3217">  
<f name="MasculineGender" xml:id="Masc" xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-3312">  
  <string>m</string>  
</f>  
</fs>  
<fs type="Article" xml:id="Article" xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-1892">  
<f name="definiteArticle" xml:id="A-def" xmlns:dcr="http://www.isocat.org/ns/dcr"  
dcr:datcat="http://www.isocat.org/datcat/DC-1430">  
  <string>A</string>  
</f>  
</fs>
```

\*A\* HK 230. f230^#1382.1 = falh0922\_pir^#9.1  
\*HL\* **Falb:1**  
\*QU\* Linz Wagner  
\*QN\* {5.2a41} Linz:nHausrv.:OÖ \*^@ FbB.WAGNER. (19xx)  
[SFb./EFb.] \*O\* Linz OÖ  
===  
\*NR\* 53E8: der Falb(e)/Falch(e) als TierN etc.  
\*LT1\* **da fâib [m+A]**  
\*LT2\* **d' fâim [pl+A]**  
\*BD/LT1\* der Falb (gelblichgraue Pferde)

**\*LT1\* [m+A] = Masculine + Definite Article**

**\*LT2\* [pl+A] = Plural + Definite Article**

(Note: gender specified in first \*LT\* applies to all \*LT\* for given entry)

# Key TEI Documents: Feature Structures for Concepts

- Can contain all and any uri for a given concept in separate <f>'s under the same <fs>;
- Reusable and extensible w/in our project and institution;
- Usable for both the dialectal data and for certain parts of the original data such as the field questionnaires that were used to elicit lexical data from speakers;
- (where possible) concepts will correspond to the sense of a given lexical item;
- LOD/RDF compatible;



# Key TEI Documents: Feature Structures for Concepts

`<fs @xml:id @type @feats>`

`<f @name="name">` (1..n)

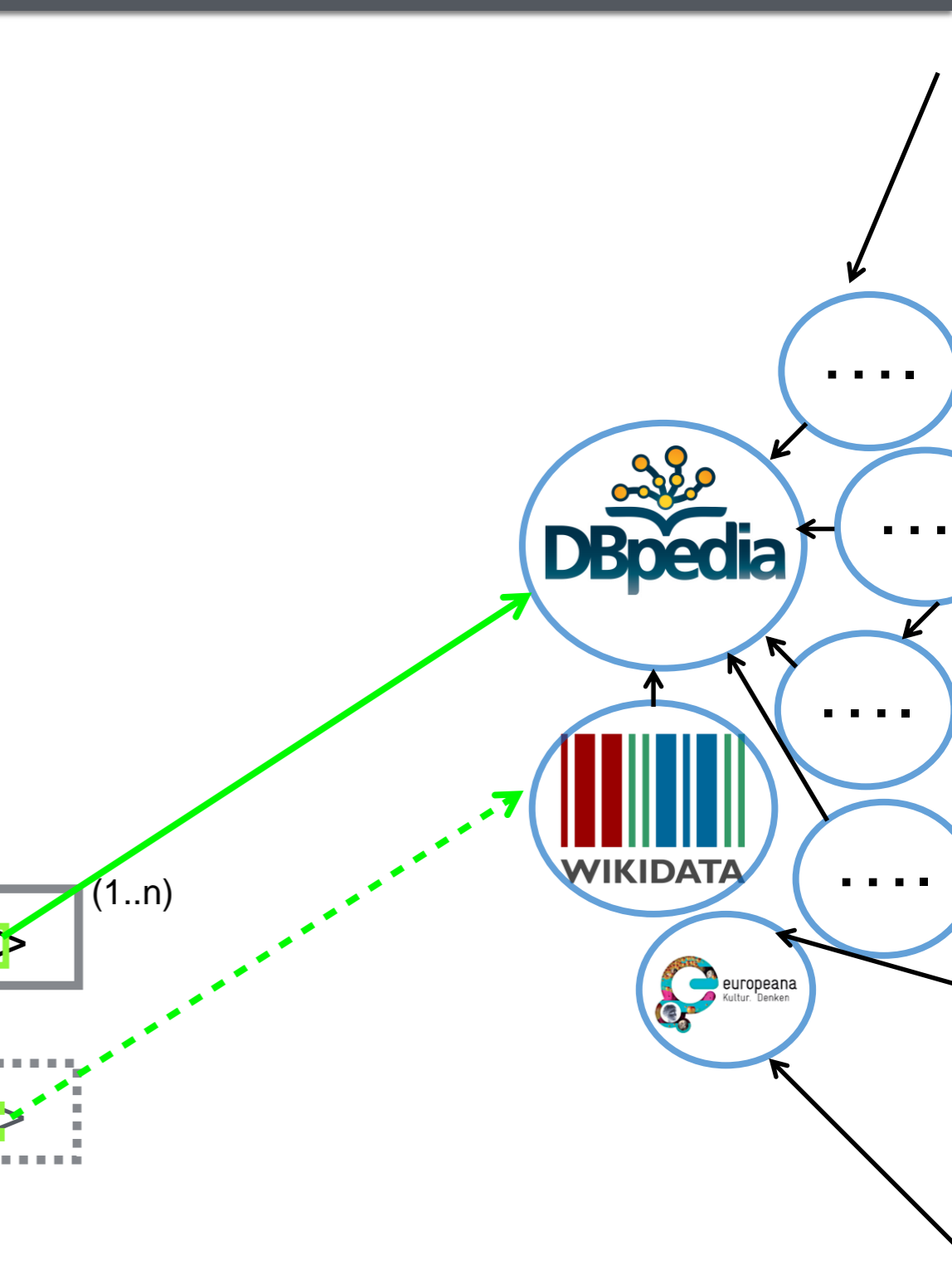
`<string @xml:lang>` (1..n)

`<f @name="entityType">` (0..n)

`<string @xml:lang>` (1..n)

`<f @name="conceptURL" @xml:id @corresp>` (1..n)

`<f @name="conceptURL" @xml:id @corresp>`



# Output with Links: Plantnames

Sample TEI record (from SQL data):

*Taraxacum officinale* Weber



```
<list type="wordforms" corresp="concept:d1e31">
...
<item>
  <cit type="example" xml:id="o151503">
    <quote>Hasenohrwaschel</quote>
    <usg type="geo">
      <placeName type="region" ref="place:Zillert">Zillert.</placeName>
    </usg>
  </cit>
</item>
...
```

## Conceptual Reference <fs>

```
<fs type="concept" xml:id="d1e31">
  <f name="name">
    <string xml:lang="la">Taraxacum officinale Weber</string>
  </f>
  <f name="name">
    <string xml:lang="de">Wiesen-Löwenzahn</string>
  </f>
  <f name="entityType">
    <string xml:lang="en">plant</string>
  </f>
  <f name="conceptURL"
    xml:id="d1e31_cURL1"
    corresp="http://dbpedia.org/resource/Taraxacum_officinale"/>
  <f name="conceptURL"
    xml:id="d1e31_cURL2"
    corresp="http://openup.nhm-wien.ac.at/commonNames/references/
scientificName/3926"/>
</fs>
```

## Place Reference <listPlace>

```
<place xml:id="Zillert" type="region" corresp="http://sws.geonames.org/2760567">
  <placeName role="main">
    <placeName full="abb" xml:lang="de">Zillert.</placeName>
    <placeName full="yes" xml:lang="de">Zillertal</placeName>
  </placeName>
  <location>
    <geo corresp="#gis_region_id-771"/>
  </location>
  <note>häufige Abk. "Zill.": Tal der Ziller samt einmündender Alpentäler im SO vom —:
#mNTir.# --- vom Sammler Michael Juffinger mit der Ziffer IV bezeichnet ---</note>
</place>
```

# Linking of dialectal data with conceptual infrastructures

## Dialectal entries: (from TUSTEP)

dialect-form: **gäb**

dialect-form: **gölj**bw****

de: **gelb**

de: **gelb**

etymological-form: **gëlb**

etymological-form: **gëlb**

dialect-place: **Kufstn.NTir.**

dialect-place: **Kl.Feistritz**

Concept:

de: **gelb**

(en: *yellow*)

## Feature structure

```
<fs type="concept" xml:id="c1e3186">  
  <f name="name">  
    <string xml:lang="de">Gelb</string>  
  </f>  
  <f name="entityType">  
    <string>color</string>  
  </f>  
  <f name="conceptURL" xml:id="c1e3186_cURL1"  
    corresp="http://dbpedia.org/resource/yellow"/>  
</fs>
```

## BabelNet entry

bn:00081866n • NOUN • Concept • Categories: Farbproduktion, Farbname

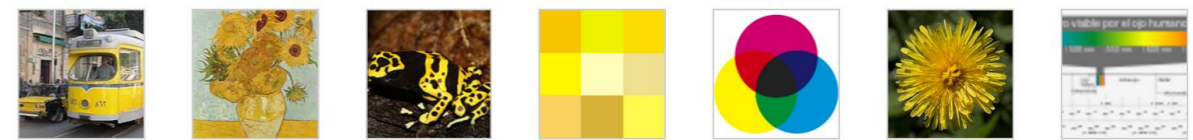
**Gelb** • Bronzegelb • Eurogelb • Goldgelb (Farbe) • Hellgelb

Gelb ist die **Farbe**, die wahrgenommen wird, wenn **Licht** mit einer spektralen Verteilung ins **Auge** fällt, bei der **Wellenlängen** zwischen 565 und 575 nm dominieren.

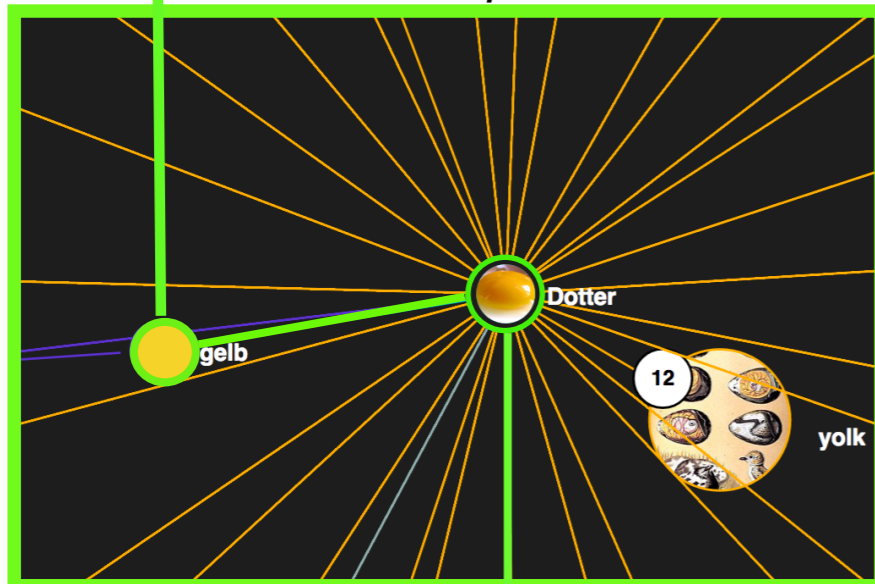
More definitions

IS A: Spektralfarbe • Grundfarbe • Farbe  
COLOR SPACE: CMYK-Farbmodell  
SOURCE: Webfarben

EXPLORE NETWORK



related concept



# Next Steps:

- Continue refining conversion process
- Integrate concepts <fs> with more output (NLP, manual)
- Begin converting groups of vocabulary
- Develop phonetic <fs> inventory for all transcription systems used in DBOE
- TEI Dictionary (???)
  - Enhanced markup for
  - Etymology; LOD semantics, inflection,...
- Conversion to RDF(?)
- Analyses of linguistic content
- Visualize data from DBOE (TUSTEP) and (DBO@ema) by geographic, temporal and lexical features

# Danke!

 austrian  
centre for  
digital  
humanities



ÖSTERREICHISCHE  
AKADEMIE DER  
WISSENSCHAFTEN

## Contact

**Jack.Bowers@oeaw.ac.at**

**<http://acdh.oeaw.ac.at>**