



# TEI Modelling of the Lexicographic Data

Krzysztof Nowak

Dorota Mika

Wojciech Łukasik



Rzeczpospolita  
Polska

Unia Europejska  
Europejski Fundusz  
Rozwoju Regionalnego



## The Dariah.Lab Project Lexicographic Resources

### The workflow

### Discussion

- delimiting lexical units
- defining senses
- representing language use
- placing in time and space

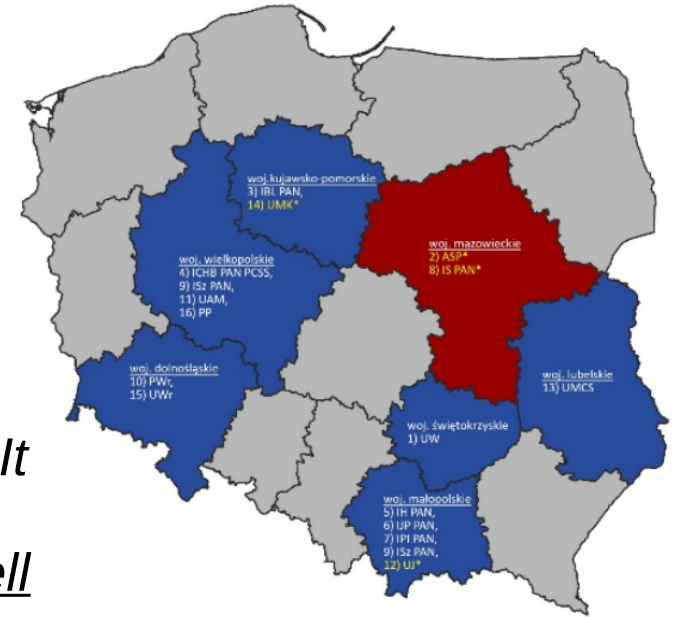


WWW: <https://lab.dariah.pl/>

Period: January 2021 - December 2023

Funding: European Union

**Goal:** *Dariah.lab is a research infrastructure for the arts and humanities built in the DARIAH-PL project. It is designed to acquire, store and integrate cultural data from the humanities and social sciences, as well as to process, visualize, and share digital resources.*  
(Source: Project's website)

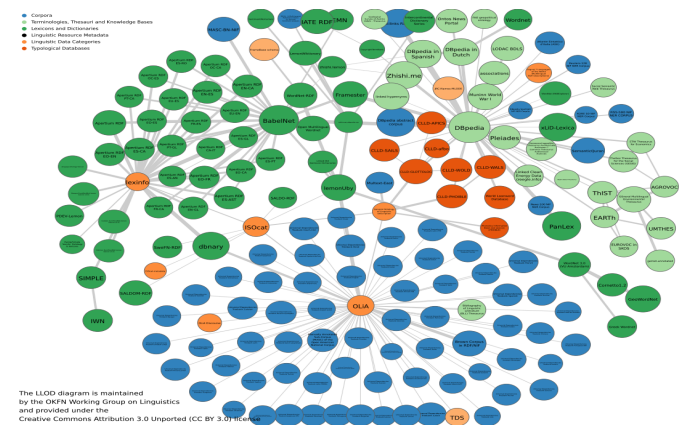


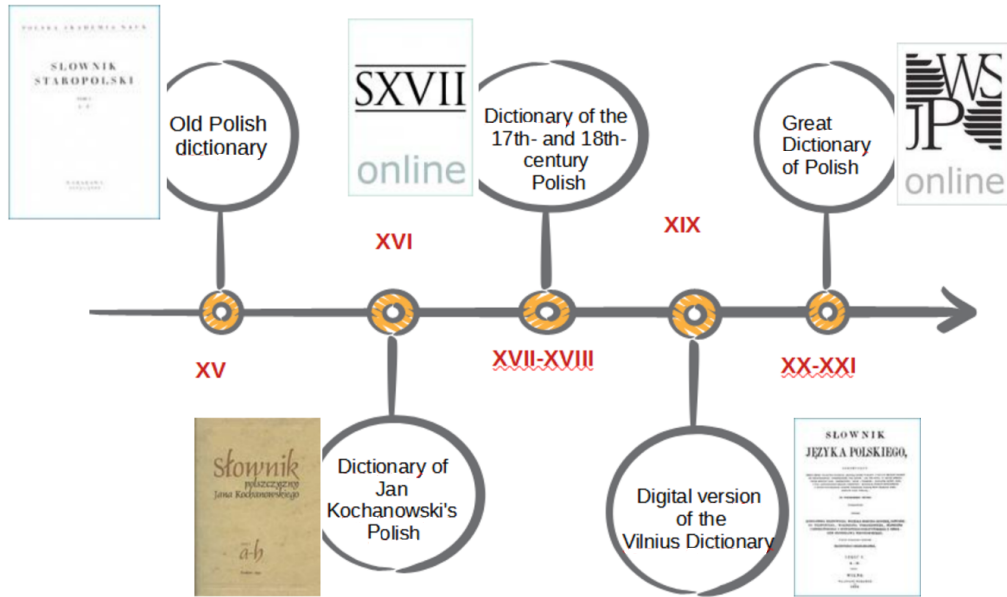
1. The Spoken Component of the Polish Corpus of 2011-2020 (in collaboration with the Institute of Computer Science PAS )



2. Integrated Access to the Institute's Dictionaries and Lexicographic Data

- web platform for data discovery
- API for automated access





## Wojciech Łukasik

- computational literary studies

## Dorota Mika

- linguistics and lexicography of Old Polish

## Krzysztof Nowak

- linguistics and lexicography of Medieval Latin

1. **BUC** rz 1. 'ktoś o wydatnych policzkach, pucłowaty, zwłaszcza dziecko': Od nŷy, fšiske byli take buce. Fšiske jejŷy ŷeći Bobowo st-gdań.

2. 'ktoś małomówny': Ostrowce bus.

3. *przez* 'ktoś niesympatyczny, za rozumiały': To taki buc, pšeiŷy i ŷy nŷy uodezvŷy Bobrka kroś.

4. 'ktoś niezaradny, niedoŷega': Tŷ bucu stari. Ale to beł buc. Ta so vzała ale buca pn i pd ok Źarnowca wej S I 81.

5. 'przez nadawane mieszkańcom wsi': Buc, bamber — tak nazywają ludzi ze wsi Oborniki.

```
<entryFree n="BUC1" type="hom">
  <form>
    <orth>BUC</orth>
    <gen>rz</gen>
  </form>
  <sense><lbl>1.</lbl><def>ktoś o wydatny
    dziecko</def>: <quote>Od nŷy/ fšis
    <usg type="geo">Bobowo st-gdań.</usg>
  <sense><lbl>2.</lbl><def>ktoś małomówny
  <sense><lbl>3.</lbl><def>ktoś
    taki buc, pšeiŷy i ŷy nŷy uodezvŷy
    kroś.</usg></sense>
</entryFree>
```

**BUC** rz 1. 'ktoś o wydatnych policzkach, pucłowaty, zwłaszcza dziecko': Od nŷy/ fšiske byli take buce. Fšiske jejŷy ŷeći Bobowo st-gdań.  
 2. 'ktoś małomówny': Ostrowce bus.  
 3. *przez* 'ktoś niesympatyczny, za rozumiały': To taki buc, pšeiŷy i ŷy nŷy uodezvŷy Bobrka kroś.

**Jałowica, Jełowica** *formy: n. sg. jałowica* 35  
 1400 *StPPP* VIII 929, XV *p. post. PF* III 286;  
 jełowica 1453 *AGZ* XIV 388; ~ *d. sg. jałowicy*  
*BZ* Deut 21, 4; ~ *ac. sg. jałowicę* 1423 *StPPP*  
 II nr 1910, 1439 *AGZ* XI 158, 1440 *AGZ* XI  
 165, etc.; ~ *i. sg. jałowicą* *BZ* Deut 21, 6;  
 ~ *n. du. jałowicy* 1446 *AGZ* XI 272; ~ *ac. du.*  
*jałowicy* 1453 *AGZ* XI 386; ~ *ac. pl. jałowice*  
 1477 *StPPP* II nr 4194, 1484 *AGZ* XV 531.

**Znaczenia: 1.**  
*iuvenca, vitula*: Pro  
 1400 *StPPP* VIII 92  
 dare debet 1423  
 terciam penam me  
 1439 *AGZ* XI 158;  
 plenam... facultate  
 iałowyczą 1440 *ib.*  
 luet iałowicza 1442  
 czelne et due iałow  
 45 Due iałowyczy trzecie  
 nyestany iałowiczą

**Jałowiec** *bot. 'jałowiec pospolity, krzew lub*  
*jego owoc, Juniperus communis L.*: *Yalowyecz*  
*bacce iuniperi* 1419 *Rost* nr 5069; *Ialowyecz*  
*iuniperus* 1437 *ib.* nr 10836, *sim.* 1460 *ib.*  
 nr 3602, *ca* 1465 *ib.* nr 3929, etc., XV *p. post.*  
*R* LIII 66, *ca* 1500 *Erz* 23; *Yalouiecz* *iunipe-*  
*rum* 1444 *R* XXIII 302, *sim. ca* 1500 *Erz* 23;  
*Jałowczą* (*GIWroc* 85v. 133v, *R* XLVII 353;  
*yalowecz*) *proicitque se et dormiuit sub vmbra*  
*iuniperii* (III *Reg* 19, 5) XV *p. pr. SKJ* I 303;  
*Jałowecz* XV *med. PulKras* 33; *Iuniperi* *vlg.*  
*jałowycz* 1456 *R* XXXIII 183; *Jałowycz* *iuniperi-*  
*um* 1463 *PF* V 12; *Yalowiecz* *iuniperius*  
 1464 *Rost* nr 4819; *Yalowecz* *iuniperus* 1472 *ib.*  
 nr 93, *sim.* 1491 *ib.* nr 11065; *Yalovyecz* *iuni-*  
*perius* 1478 *ib.* nr 2110; *Yalovyecz* *ib.* nr 2266;  
 \**Yalovyecz* *iuniperus* *ca* 1500 *ib.* nr 2060.

**Jałowię** *'cielę, cielątko, vitulus'*: *Vitulum* *al.*  
*jałową* 1429 *ArchCastrCrac* III 188; *Petrus...*  
 55 *tres vaccas, duo jałowata presentis anni* *al.*

**BAZGRAĆ (SIE)** *Formy: 3 os lp*  
*czter* *baźgro* || *baźgże* ||- *Mszana G* *lim*;  
 ~ *baźglać* ||- *Mszana G* *lim*; ~ *bez-*  
*grać* ||- *ok Sejn suw PKJP* II 2 s 51;  
 ~ *baźgżoś* ||- *ok Sejn suw PKJP* II 2 s 51;

**CHŁOPYSIO, CHŁOPYSIU, CHŁO-**  
**PYŚ** *Formy: Chłopyśo* ||- *Libiszów opocz*  
*3 os lp czter* *bejzgra* ||-  
*RŁTN* XXI 157; *Chłopyśo* *Kuj Kuj* II 285; *Chłopyśo* *JJP* II 2 s 51; ~ *pazgrać*  
 ~ *Chłopyśu* ||- *Domaniewek łącz PJPAN* ; *Łopienno wąg* *PKJ*  
*XXXVII* 309; ~ *Chłopyś* ||- *Libiszów opocz*  
*RŁTN* XXI 157; ||- *Domaniewek łącz PJPAN*  
*XXXVII* 309; *Kramsk koniń; Chłopyś Pod-*  
*różna złotow; ~ W lp chłopyśu* *Peary łącz*  
*LL* VI 1-2 s 85.

**Znaczenie: 'chłopczyk':** *Chłopyś*  
*myśl'i, że ty tu som doż rade* *K*  
*Chłopyśu, pūy, dwoęstańeš japk*  
*złotow; Libiszów opocz RŁTN*  
*maniewek łącz PJPAN XXXV*  
*łącz LL* VI 1-2 s 85; *Kuj Kuj* I

**I. BOSAK, BOSAKA, BÓSAKI** *przy-*  
*słów* *Formy: Bosak* [*rzd*] *Krościenko*  
*n-tar* *MAGP* VIII m 364 s 56; *bosok*  
*Trojanów miń-maz jw* m 364; ~ *bosoka*  
*Domaniewek łącz PJPAN XXXIII* 47;  
*Domaszew gar; Złotniki kal; bōsoka*  
*Obra wolsz; bosaka* [*rzd*] *Bąków rad*  
*MAGP* VIII m 365; ~ *bosaki* *Włynice*  
*radomsz RŁTN* XX 249; ~ *typ* *bossaka:*  
*Bąków rad* *MAGP* VIII m 364 s 56; || *bo-*  
*saka jw* m 365 s 58; *Mokrzysz częs* *MAGP*  
*VIII* m 364; *Bolmin kiel RŁTN* XX 250;  
*Kajetanów ilż jw; ||- Jelnia opocz jw; ~*  
*typ* *boscaka: Alojzów ilż* *MAGP* VIII  
 m 365; || *bosaka* *AWK* III 377; *Opatko-*



## TIFF

Boruta 395

stawiają naczyńia z potrawami Przyj-  
Ludu VI 127. JO  
II. BORUTA 'nazwa krowy': Boruta  
MaguŃy i ok [Grodno ZSRR] Polgów II 169.  
JO  
Boruwać zob. I. BOROWAĆ  
Borwinek zob. BARWINEK  
BORYJÓWICA Forma: Boryjuwica  
Przewrotne rzesz.  
Znaczenie: 'gnojówka'; jw. JO  
Boryjówica zob. BORYJÓWICA  
BORYJÓWKA 'gnojówka': Boryjówka  
Przewrotne rzesz. JO  
BORYKAĆ (SIE) 'plądrować, szuka-  
jąc przetracać': Borykać Krak K I 108.  
Borykać się I. 'szamotać się, walczyć; sięć  
w zapasy': Borykać się Krak jw.  
2. 'pokonywać trudności, zmagać się z  
przeciwnościami; męczyć się': Tagęz  
śe borykou aż do roku tyśiune żewensot  
štyrcyot pięć Niebory [Cieszyn Csz] PJP  
I 46; Zyl z pracy rąk swoich, a I ty  
pracy nié było dzie pójnialé i tak się w  
ty bydzicie swoji borykál Żurawicki  
przew LL I 3 s 26; Cnye żyće cówek  
śe tylko borykáj, nié śe użyz Samocieć  
dq̄b-tar; Se żużok cnye żyće borykou  
ji tero na stare lata nié jé ni mo  
Domaniewek łącz PJPAN XXXIII 47;  
Krak K I 108; Złotniki kal: ~ B. się z  
czym: Boryko śe z rzuŃnymi nieścycamy  
Študzienice rad RZTY XX 249.  
3. 'starać się, zabiegać o coś': Juz  
wypowegali nié máz-roboty, a uto  
robote śe borykałi luže, byli śe  
Niedzica n-tar ZYUJ CCLXXVIII 126.  
JO  
BORYKI bjp 'zapasy': Pastuch [...] máz  
w boryki z kudiatym? ['wiłoka-kiem']  
BodzentyŃ kiel Lud I 204. JO  
BORYN \*starcie do przewożenia  
żyłczy z hali do wsi": Boryna  
Dzianisz n-tar PJP IV 41. JO  
BORYNKA 'drzbanek': Borynka  
[|] borynka Rabka-Zaryte n-tar; —  
'błazana bańka na mleko, jagody':  
Borynka Za-woja wad PJP XIII 64. JO

Borzyska

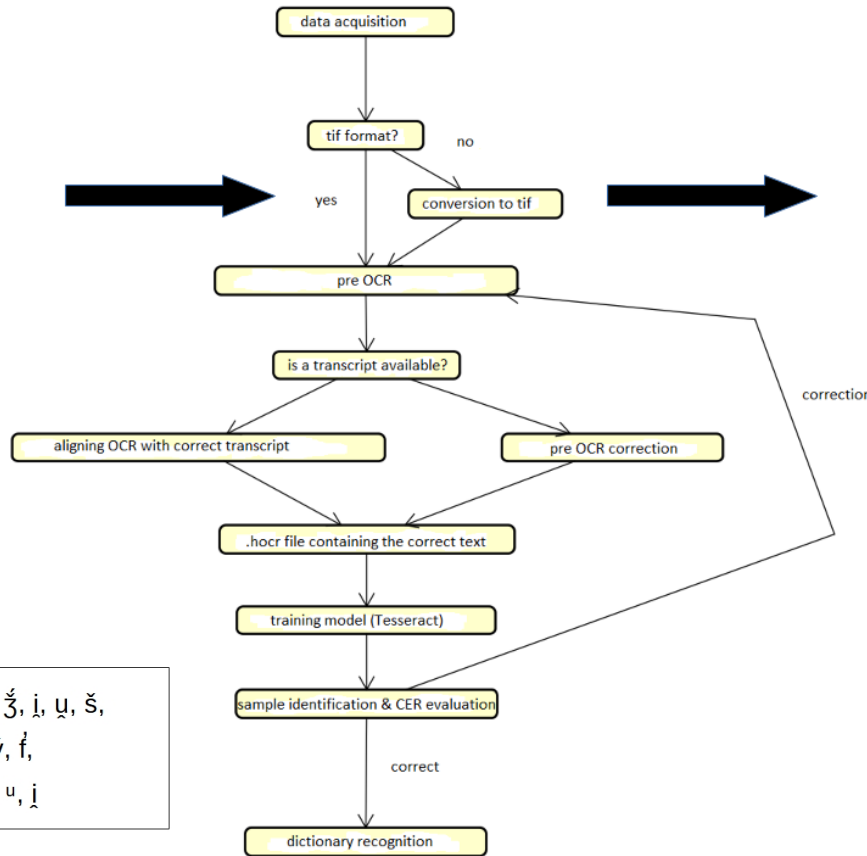
II. Borys zob. PORYS  
BORYSZ 'poczętnak po dobligum  
targu': Jak zabił na ten borys, tu pój  
Degucie sus; Boryša mus'leli my złaŃfo  
wypic, n'awet żyt jak kupował st'awił  
čv'artka ZopatoŃszczyzna [Wilno ZSRR]  
PJFAN LXI 121; DobryŃseno D biało;  
Ślavecjin čhozec [Michaliński, Wilno  
ZSRR] ZPSSJ IV 133.  
Por. BARYSZ JO  
BORYSZEK 'zdř od bór w zn I';  
Boryšek Kubeczi ravoie PKJP III 3  
s 150. JO  
Borzyska zob. I. BORZYSKA  
BORYSZOWAĆ 'hulać': Jak się dzie-  
dzie dowiedził, zo rzońca tak borysuje,  
dał mu trzy dni folgi Turbio radz-podł  
PF VI 219. JO  
BORYTASY cwykto w Im 'rodzaj  
ornaments zdobięcego spodnie góralskie';  
Męczyzjni nosili na terenie całego Spisza  
spodnie z bialego sukna [...] zdo-  
bione haftowanymi porzecziami [...] lub [...]  
uwtymi z dwubarwnego sznur [...] albo z  
dwubarwnych tasemek z dodatkim  
haftu zw. „borytasami” N-tar Star-Przew  
34; N-tar pd-wok ASL V I s 12. JO  
Boryna zob. I. BURZYNA  
BORZAN 'mieszaniec lam': Bożańe  
a leścio zo jo samo Kasz S I 59. JO

combining characters: 0, 3, 3̇, 3̈, 3̉, j, u, š, č, ž, ś, ǰ, k, g, h, p, b, t, d, à, â, v, f, w, m, n, r, t, n, d, o, â, y, â, j, õ, u, j

## OCR

Boruta 395

stawiają naczyńia z potrawami Przyj-  
Ludu VI 127. JO  
II. BORUTA "nazwa krowy": Boruta  
MaguŃy i ok [Grodno ZSRR] Polgów II 169. JO  
Boruwać zob. I. BOROWAĆ  
Borwinek zob. BARWINEK  
BORYJÓWICA Forma: Boryjuwica  
Przewrotne rzesz.  
Znaczenie: "gnojówka"; jw.  
Boryjówica zob. BORYJÓWICA  
BORYJÓWKA 'gnojówka': Boryjówka  
Przewrotne rzesz. JO  
BORYKAĆ (SIE) "plądrować, szuka-  
jąc przetracać": Borykać Krak K I 108.  
Borykać się I. 'szamotać się,  
walczyć; sięć w zapasy': Borykać się  
Krak jw.  
2. "pokonywać trudności, zmagać się z  
przeciwnościami; męczyć się": Tagęz  
śe borykou aż do roku tyśiune  
żewensot štyrcyot pięć Niebory  
[Cieszyn Csz] PJP I 46; Zyl z pracy rąk  
swoich, a I ty pracy nié było dzie  
pójnialé i tak się w ty bydzicie swoji  
borykál Żurawicki przew LL I 3 s 26;  
Cnye żyće cówek śe tylko borykou, nié  
śe użyz Samocieć dq̄b-tar; Se żużok  
cnye żyće borykou ji tero na stare lata  
nié jé ni mo Domaniewek łącz PJPAN  
AXXXIII 47; Krak K I 108; Złotniki kal:  
— B. się z czym: Boryko śe z rzuŃnymi  
nieścycamy Študzienice rad RLTN XX  
249.  
3. 'starać się, zabiegać o coś': Juz  
wypowegali nié máz-roboty, a Io robote  
śe borykałi i luze, byli śe Niedzica  
n-tar ZNUJ COLXXVIII 120. JO  
BORYKI bjp "zapasy": Pastuch [...] máz  
w boryki z kudiatym? [wilkoła-kiem]  
BodzentyŃ kiel Lud I 204. JO  
BORYNA "starcie do przewożenia  
żyłczy z hali do wsi": Boryna  
Dzianisz n-tar PJP IV 41. JO  
BORYNKA "drzbanek": Borynka  
[|] borynka Rabka-Zaryte n-tar; —  
'błazna bańka na mleko, jagody':  
Borynka Za-woja wad PJP XIII 64. JO





## Experiments

- grobid-dictionaries
- other ML solutions

## Solution: custom Python (\*CLARIN-PL) and XSLT scripts

- punctuation marks
- anchoring strings: labels, abbreviations, initials etc.
- location
- no font style info

## Motivation

- (slightly more) formal resource's description
- documenting
  - annotation choices and semantics
  - original rendering
- generating scheme and documentation

## Principles

- preserving textuality for quoting and display
- exposing data to enrichment
  - aligning dictionary parts
  - dates, places, sources etc.
  - reconciling meta-languages

```

<schemaSpec ident="ijpsgp"
             source="ijp_all.compiled.odd">
  <elementSpec ident="entryFree" mode="change">
    <gloss xml:lang="pl">hasło</gloss>
    <desc xml:lang="pl">hasło słownikowe</desc>
    <attList org="group">
      <attDef ident="type" mode="change">
        <valList type="closed" mode="change">
          <valltem ident="main" mode="change">
            <gloss xml:lang="pl">zwykle</gloss>
            <desc xml:lang="pl">zwykle hasło</desc>
          </valltem>
        </valList>
      </attDef>
    </attList>
  </elementSpec>
</schemaSpec>

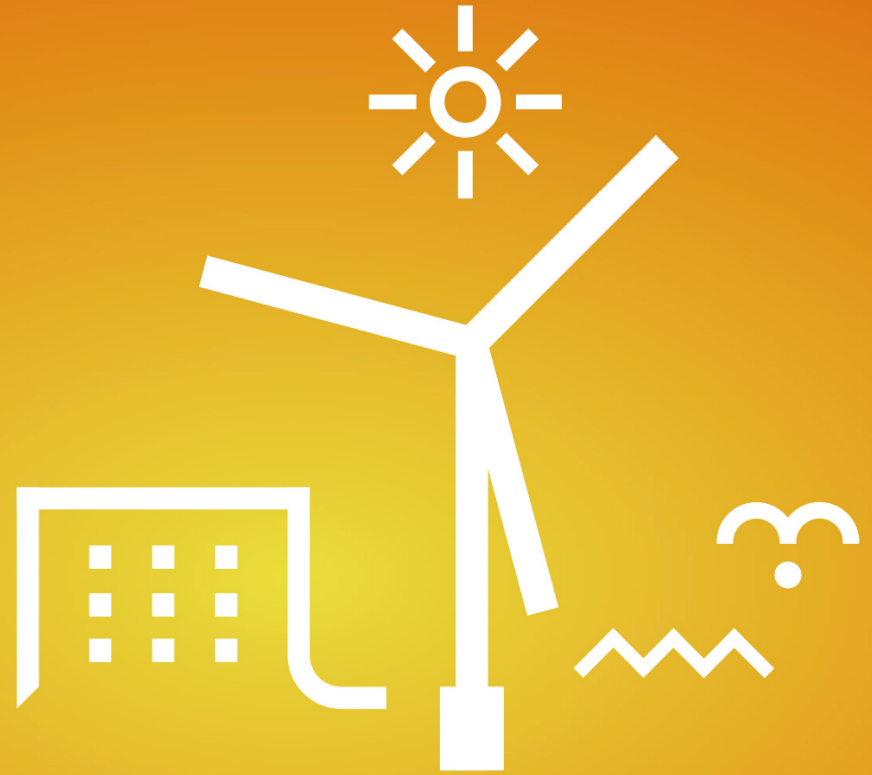
```

**Delimiting lexical units**

**Linguistic evidence:  
time and space**

**Taking stance: critical  
apparatus**

**Citation strategies**



\*BAŁAJKA

```
<entryFree type="empty">  
  <form type="lemma">  
    <lbl>*</lbl>  
    <orth norm="bałajka">BAŁAJKA</orth>  
  </form>  
</entryFree>
```



'a balalaika'

```

<!-- main entry: a participle -->
<entryFree>
  <!-- related entry: a derived noun-->
  <re/>
  <!-- related entry: a derived noun -->
  <re/>
  <!-- related entry: a derived prepositional phrase -->
  <re/>
</entryFree>

```

**GIĘTY 1.** *imb od cz giąć (zob.)*

**2.** *'sprężysty, uginający się, elastyczny; giętki':* *b̄ył n̄e zūom̄ie śe, j̄e ḡante chojn 5.*

**Gięta** *w użyciu rz, dziec 'o łodzie uginającym się przy chodzeniu po nim':* *n̄e l̄at̄aj̄va j̄us p̄o ḡeće, b̄o j̄esce χtury z n̄az uok̄a-p̄e śe v v̄oż̄e d̄qb-tar 5.*

**Gięte** *w użyciu rz 'nosidła':* *tarn PAE IV m 241; boch jw.*

**Po giętu** *'w pozycji zgiętej, pochylonej':* *to my tag\_muśeli χ̄ożīć p̄o\_ḡentu, żeby śe n̄e zaduśīć ciesz ZNUJ 495 s 113. AN*

```
<sense>
  <lbl>~</lbl>
  <usg type="colloc">jałowcowe drzewo</usg>
  <usg type="dom">bot.</usg>
  <def xml:lang="pl">jałowiec pospolity</def>
  <def xml:lang="la">Juniperus communis L.</def>
</sense>
```



'juniper tree'

```
<sense xml:id="Sstp_Jałmużnik_sense_1">
  <def xml:lang="pl">dający jałmużnę lub z urzędu
    rozdzielający jałmużnę</def>
  <def xml:lang="la">qui eleemosynam dat vel ex officio
    distribuit</def>
  <lbl>(?)</lbl>
  <certainty locus="value" degree="0.5" match="orth"/>
</sense>
```

... word senses

```
<form type="lemma">
  <orth norm="grzycka">GRZYCKA</orth>
  <lbl>[?]</lbl>
  <certainty degree="0.5" locus="value" match="orth"/>
</form>
```

... word forms



```

<cit xml:lang="pl">
  <quote>Boryna</quote>
</cit>
<usg type="geo">
  <placeName>
    <settlement>Dzianisz</settlement>
    <lbl type="geo">
      <region type="county" when="1986">
        n-tar
      </region>
    </lbl>
  </placeName>
</usg>
<cit>
  <bibl>
    <title type="abbr">PIJP</title>
    <biblScope unit="volume">IV</biblScope>
    <biblScope unit="page">41</biblScope>
  </bibl>
</cit>

```



Coordinates: 49°19'53"N 19°52'2"E

<b>Country</b>	Poland
<b>Voivodeship</b>	Lesser Poland
<b>County</b>	Tatra
<b>Gmina</b>	Kościelisko

```

<cit>
  <quote xml:id="quote-ialmuzna-1"> ...</quote>
  <bibl>
    <title type="abbr" xml:id="BZ-1">BZ</title>
    <biblScope unit="chapter">2</biblScope>,
    <biblScope unit="line">16</biblScope>
  </bibl>
</cit>

```

```

<cit>
  <ref corresp="#quote-ialmuzna-1" type="sim">sim.</ref>
  <bibl>
    <ref target="#BZ-1">ib.</ref>
    <biblScope unit="chapter">12</biblScope>,
    <biblScope unit="line">9</biblScope>
  </bibl>
</cit>

```

```

<cit>
  <quote>Tegodlya,
    <seg xml:id="seg-4">kyedy czynysch yamvzna</seg>
    <note>
      (<gloss target="#seg-4" resp="#lexicographer">
        cum... facis eleemosynam
      </gloss>
      <bibl type="bible">
        <title>Mat</title>
        <biblScope unit="chapter">6</biblScope>
        <biblScope unit="line">2</biblScope>
      </bibl>)
    </note>, nye day przed sobą trąbycz ...
  </quote>
  <bibl>
    <title type="abbr">Rozm</title>
    <biblScope unit="page" from="271" to="272">271-2</biblScope>
  </bibl>
</cit>

```

```
<cit>
  <quote>
    Ialmvsznym
    <corr resp="#lexicograph" xml:id="corr-1">nacznem</corr>
    <note>(
      <lbl>leg.</lbl>
      <sic corresp="#corr-1">nędznem</sic>
    )</note> nye daval
  </quote>
</cit>
```

## Keep track of advances in token-level text structuring

- sharing tools and expertise

## Discuss with content creators

## Take annotators' feedback into account

## Work with ODD

- describe the resource
- document annotation choices
- explain assumptions about document semantics
- generate multi-lingual documentation

## Know when to stop