



# The Polish Parliamentary Corpus

Maciej Ogrodniczuk | Linguistic Engineering Group  
Institute of Computer Science  
Polish Academy of Sciences

CLARIN-PLUS Worskhop on Working with Parliamentary Records  
Sofia, 27–29 March 2017

# The Polish Parliamentary Corpus

## In a nutshell:

- a large collection of linguistically annotated documents from the proceedings of Polish Parliament, Sejm and Senate
- based on the Polish Sejm Corpus prepared in October 2011 and recently updated by CLARIN-PL
- automatically created multi-layer annotation compatible with the National Corpus of Polish
- 42.4 GB compressed data, 300M segments
- available at: <http://clip.ipipan.waw.pl/PPC>

# Corpus data

## Data sources:

- 1993–: editable PDFs, other formats available (internally)
- 1989–93: scanned paper transcripts (OCR-ed and manually verified)

## Corpus format and structure:

- stand-off TEI P5 annotation
- generated with in-house tools:
  - Morfeusz SGJP (text structure, utterance-level segmentation, tokenization, lemmatization)
  - Pantera (disambiguated morphosyntactic description)
  - Spejd (syntactic words and syntactic groups)
  - Nerf (named entities).

# Source scans

(Początek posiedzenia o godz. 17 min. 05)

(Na solę uchodzącą ustnie oszacujemy okłamań delegacje Związku Socjalistycznych Republik Radzieckich, Czechosłowackiej Republiki Socjalistycznej oraz Niemieckiej Republiki Demokratycznej)

(Przewodniczący na posiedzeniu Marszałek Sejmu Czesław Węchoł)

**Marszałek:**

Otwieram posiedzenie.

Na sekretarzy delegacji państw Reginę Pawlikowską i Zdzisława Kurawskiego.

Protokół i listę mówców prowadzi poseł Regina Pawlikowska.

Przez Posłów Sekretarzy o zajęcie miejsca przy stole prezydialnym.

Przewodniczący Sejmu uprasza wierzności o obecność na posiedzeniu państw, których naszawka będą umieszczone w zaliczniku do protokołu dzisiejszego posiedzenia.

Ustalony przez Przewodniczącego Sejm porządek dzienny obecnego posiedzenia obejmuje Uroczyste uświetnienie 20-lecia Polskiej Rzeczypospolitej Ludowej.

Kier w sprawie porządku dziennego są jakiegdyż uwagi?

Nie słyszę.

Uważam, że Sejm porządek dzienny, przedstawiony przez Przewodniczącego Sejmu, zatwierdził.

Obywatelsi Państwie! Uroczyste nasze posiedzenie odbywa się w szczególności podniosłym momencie — w przeddzień 23 rocznicy powstania Polski Ludowej. W posiadaniu naszym bierzmy udział dostojnie i diermy nam Goście. Niech mi wspaniałe wolno, otwierając posiedzenie, w imieniu Was wszystkich, Obywatelsi Państwie, pod Ich adresem przede wszystkim skierować moje słowa.

Witam serdecznie przybyłą na nasze posiedzenie delegację Związku Socjalistycznych Republik Radzieckich w osobach Sekretarza Generalnego Komitetu Centralnego Komint-

stycznej Partii Związku Radzieckiego Towarzysza Leonida Brzniewa i Przewodniczącego Prezydium Rady Najwyższej ZSRR Towarzysza Nikołaja Podgornego. (Oszczepcie okłamań)

Witam serdecznie delegację Czechosłowackiej Republiki Socjalistycznej w osobach Prezydenta Czechosłowackiej Republiki Socjalistycznej Towarzysza Ladka Swobody i Sekretarza Komitetu Centralnego Komunistycznej Partii Czechosłowacji Towarzysza Gustawa Humka. (Oszczepcie okłamań)

Witam serdecznie delegację Niemieckiej Republiki Demokratycznej w osobach Prezosa Rady Ministrów i następcy Przewodniczącego Rady Państwa Niemieckiej Republiki Demokratycznej Towarzysza Willi Stophę oraz członka Biura Politycznego Sekretarza Komitetu Centralnego Niemieckiej Socjalistycznej Partii Jedności Towarzysza Ericha Honckera. (Oszczepcie okłamań)

Witam wszystkich obecnych.

Przystępujemy do porządku dziennego.

Przez o zabranie głosu i Sekretarza Komitetu Centralnego Polskiej Zjednoczonej Partii Robotniczej pości Władysława Gomułka. (Dziękuję okłamań)

**I Sekr. KC PZPR**

**poseł Władysław Gomułka:**

Obywatelsi Państwie! Wysoka Iabó Prząd 23 Iyłu otwiera się nowa karza dziejów Polski.

W Chelmie, w pierwszym uwolnionym od okupanta mieście polskim powolany został z woli Krajowej Rady Narodowej — podziemnego parlamentu polskiej Narodowej — Polski Komitet Wyzwolenia Narodowego. Jego siedziba i szereg pierwszą stoliczną kolebką Polski Ludowej było miasto Lublin.

Manifest Polkiego Komitetu Wyzwolenia Narodowego obwieszczały odródnienie Polski stał się zapowiedzią nowej epoki w historii naszego narodu.

7

I sesja Sejmu — 4 października z 21 lipca 1969 r.  
Uroczyste uświetnienie 20-lecia PRL.

8

**I Sekr. KC PZPR poseł Gomułka**

Obóz demokracji polskiej, którego przewodnicząca Polska Robotnicza włożył w słowa Manifestu Lipcowego całe swe historyczne doświadczenie, wyrabił w nich wieloletnie działania i pragnienia ludu pracującego. Program Manifestu stał się fundamentem, na którym krapka po kropce, pietra za piętrem, od 23 Iyłu wzniesiony gmach nowej Polski — Polski socjalistycznej.

W tych pierwszych dniach wolności, gdy lud pracujący rozpoczynał budowę swojej państwowości, nasz ówczesny Kład — Polski Komitet Wyzwolenia Narodowego — tyle mił ziemi pod stopami, ile jej co dzień zdobywał wywalczony korpus radziecki i dywizje polskie w krwawych bojach z hitlerowskim okupantem. Tyle mił siał, ile mu udzieliły swoje wyzwoleńce wsi i miasta, ile mu jej stał i siołko Iudzie pracy, skupiający się wokół swej wsiady.

Wówczas w Lublinie, kiedy zakładałmy społeczne, polityczne i prawne podstawy pod nową Ludową Rzeczpospolitą, zrewą najwspanialszą była wojna z najciężką hitlerowską wojną, która naród polski toczył nieustannie od I września 1939 r., i która dokoła osesa była od zakończenia. Gdy PKWN schodził za sceny historycznej, ustępując miejsca Rządowi Tymczasowemu, Iudzie w Wolskiej, Narwi i Wisłocie, 75% ziemi polskiej znajdowało się jeszcze pod władzą okupanta, 75% ludności polskiej narazem było osła na postępującą eksterminację. Dorołek materialny Polaków ulegał dalekemu niszczeniu przez okupanta. Toż Igrawo wojny, ogień i przemieszczenia wywołania i — wzmocnia udziału narodu polskiego w walce o zmiadziasta Niemiec hitlerowskiej, tradycjonalny jako naszebie za daniem nowej władzy, najwzyszy nasz dany nam przez naród.

W dniu swych narodziń ludowe państwo polskie miało 65 tysięcy żołnierzy na froncie i Armie Wojska Polkiego stworzona przez Iewice polskie na ziemi radzieckiej, przy jak-40 tysięcy młodego żołnierza w powstających dopiero formacjach, nie licząc partyzanckiego ruchu podziemnego. Po pięciu miesiącach zwycięzcznie uderwionych żołnierzy w szeregu zjednoczonego ludowego Wojska Polkiego. Pod koniec wojny, wiosną 1945 roku, milionały pod brońą prawie 400 tysięcy Iudzi, z czego 185 tysięcy wzięło udział w wielkich wiosennych bitwach, które przyniosły wolność Polsce i zadecydowały ostatecznie zwycięstwem nad hitlerowską tyranii.

Udział ludowego Wojska Polkiego w ofensywie styczniowej, a zwłaszcza w walkach na Pomorzu w Iytmu i marcu 1945 roku, wazył na przebiegu operacji wojennej na froncie wschodnim i przyczynił się w sposób istotny

do skrócenia okupacji ziem środkowej i północnej Polski. Armie polskie, uczestniczące w ofensywie kwietniowej, stanowiącej 8 do 10% ogółu wojsk walczących w gigantycznej bitwie między Siedmią Brygadą Armii Radzieckiej, która przeszła do ostatecznej klęski Hitlera. Totek fakt, iż wianie nasz żołnierz, Iudzy z sojusznikami stał obok żołnierza radzieckiego w Berlinie, stał się symbolem dzisiejszej sprawiedliwości.

Ten udział Polaków w ostatecznym zwycięstwie nad hitleryzmem odegrał wielką rolę w historycznych rozstrzygnięciach, mających na celu zabezpieczenie przytulności Polski i Europę przed zaborczością niemiecką, rozstrzygnięciach, które przywróciły Polsce jej ziemie piastowskie po Odrze i Nysie Łużycką i wyzwały Iuskytu po Szczecin.

Sprawiedliwie dla Polski rozstrzygnięcie w Poznaniu, mające u swych podstaw decyzje konferencji trzech mocarstw w Jakim, naród polski opłacił straszną daninę krwi, cierpień i strat w czasie II wojny światowej. Rachunek zbrodni hitlerowskich dokonywanych na narodziu polskim jest przerażający: 6 milionów pologkich i pomordowanych — 22% historycznej substenjacji narodu — to Iudzy, których wstrząsająca wymowa jest nie do natęrcia. A przecież nie wyczerpały one wcale agrum mił strat Polski: Z każdego Iysiska obywateli pologkich, zamieszkałych w Polsce w 1939 r., 220 zostało aszorekownie do kraja. Zwycięzcy do ziemskiego inwalidzacji; na gruzlice zapadło około 100, do osobów pracy i na przy- masowe roboty wywieziono do kraja. Zwycięzcy zabawiono całego obywateli i wysiedlono — 90, pozostaństwo bez dachu nad głową 150.

Całkowitemu zniszczeniu uległo prawie 40% Iudzi naszego Iudzi, wartość Iudzi dewaluacji przez dwóch poprzednich pokoleń Polaków.

Ten straszący bilans wyrażał jedynie częściową realizację zbrodniczych planów okupanta, planów całkowitej likwidacji narodu pologiego jako narodu zdolnego do życia. Zwycięzcy — jak zwycięca okłamał Hitler w przeddzień napisania na Polskę — nie miał odwazy się szedź.

Dziś raz jeszcze przypomniałoby, że jeśli te obędne plany pełnej eksterminacji narodu pologiego nie zostały ujęrzwionione, zawiądujemy to w stopniu decydującym w Związkuwi Radzieckim, jego Iobsterkistki armii.

Przyjęcie wyzwolenia o Iabdy tydzień omawiała dla nas uradowana Iyca, 20 tyłam- com Iudzi — tylko bowiem przedmiotem Iudziwo- wania hitlerowskiego Iudziwo- wania okupacji w naszym kraju.

Główną rolę w zwycięstwie nad faszyzmem odegrała Armia Radziecka. Ona to — jak Iudziwka Iudziwo- wania armii Iudziwo- wskiej — zadła prawie 80% wszystkich strat pologich przez hitlerowskie wojsko Iudziwo- wania w okresie od I września 1939 r. do 29 kwietnia 1945 r.

# Linguistic annotation format

## Text structure:

```
<teiCorpus>
  <xi:include href="corpus_header.xml"/>
  <TEI>
    <xi:include href="header.xml"/>
    <text>
      <body>
        <!-- ... -->
        <div xml:id="txt_7-div">
          <u xml:id="txt_7.1-u" who="#The_Speaker">
            <!-- ... -->
          </u>
          <u xml:id="txt_7.2-u" who="#MP_Jan_Nowak">
            <!-- ... -->
          </u>
        </div>
      </body>
    </text>
  </TEI>
</xi:include>
</teiCorpus>
```

# Linguistic annotation format

## Sentence- and token-level segmentation:

```
<p xml:id="segm_txt_7.1-u"  
  corresp="text_structure.xml#txt_7.1-u">  
  <s xml:id="segm_txt_7.1-u.1-s">  
    <!-- ... -->  
    <!-- Proszę -->  
    <seg xml:id="segm_txt_7.1-u.4-seg"  
      corresp="text_structure.xml  
        #string-range(txt_7.1-u,0,6)"/>  
    <!-- państwa -->  
    <seg xml:id="segm_txt_7.1-u.5-seg"  
      corresp="text_structure.xml  
        #string-range(txt_7.1-u,7,7)"/>  
    <!-- ... -->
```

# Linguistic annotation format

## Morphosyntactic annotation:

```
<seg xml:id="morph_txt_7.1-u.45-seg"
      corresp="ann_segmentation.xml
      #segm_txt_7.1-u.45-seg">
  <fs type="morph">
    <f name="orth">
      <string>państwa</string>
    </f>
    <!-- państwa [7,7] -->
    <f name="interp">
      <fs type="lex" xml:id="...">
        <f name="base">
          <string>państwo</string></f>
        <f name="ctag">
          <symbol value="subst"/></f>
        <!-- ... -->
```

# What's next?

## CLARIN-PL activities:

- processing data with newest analytic tools
- live corpus
- more annotation layers
- better search and presentation
- rich(er) metadata
- consultation of other versions
- audio/video linking
- more data!



# More data

## Which data?

- high-quality scans prepared by the Sejm library
- manually cleaned
- 2000 documents, 90K pages from 1918–1989
- two political periods most interesting for the humanities

## How to process them?

- human-assisted OCR (area and typo corrections)
- scripts for structure and error detection