



„KORBA”

**Elektroniczny korpus tekstów
polskich z XVII i XVIII w. (do 1772 r.)**

Pracownia Historii Języka Polskiego
XVII i XVIII wieku IJP PAN,
Instytut Podstaw Informatyki PAN

Podstawowe informacje o projekcie

- Projekt realizowany przez IJP PAN we współpracy z IPI PAN, finansowany w ramach Narodowego Programu Rozwoju Humanistyki na lata: 2013-2017
- Kierownik projektu: prof. dr hab. Włodzimierz Gruszczyński
- Planowana objętość korpusu to 12 mln segmentów
- Korpus historyczny rozszerzający Narodowy Korpus Języka Polskiego
- Kryptonim: KORBA (KORpus BArokowy)
- <http://korba.nlp.ipipan.waw.pl/login/?next=/> (obecnie dostępny dla członków zespołu pracującego nad korpusem oraz w Pracowni Historii Języka Polskiego XVII i XVIII w. IJP PAN w Warszawie)

Zakres znakowania

- Znakowanie socjolingwistyczne (informacje o autorze, wydawcy, tłumaczu) – umieszczone w „metryczce” każdego tekstu.
- Znakowanie stylistyczno-genologiczne – umieszczone w „metryczce” każdego tekstu.
- Znakowanie strukturalne i tekstowe, umożliwiające dokładną lokalizację cytatu w tekście (paginacja; rozdziały, części, księgi itp.; notki marginesowe, przypisy itp.; podpisy pod ilustracjami; także oznaczenie wszystkich wtrętów obcych, z podziałem na języki).
- Znakowanie morfosyntaktyczne
 - 0,5 mln segmentów – ręczne
 - reszta – automatyczne

Znakowanie socjolingwistyczne

[Powrót do strony głównej](#)

Edytuj metadane

Id: *	<input type="text" value="PetrSInst"/>
Tytuł: *	<input type="text" value="Instrukcja albo nauka jak się"/>
Autor: *	<input type="text" value="Sebastian Petrycy"/> <input type="checkbox"/> Anonimowy
Tłumacz:	<input type="text"/> <input type="checkbox"/> Anonimowy
Drukarnia:	<input type="text" value="Mikołaj Lob"/>
Rok:	<input type="text" value="1613"/>
Rok wydania niepewny (wydano nie wcześniej, niż podany rok):	<input type="checkbox"/>

Miejsce: *	<input type="text" value="Kraków"/>
Region: *	<input type="text" value="Małopolska"/>
Wydanie u Współczesne:	<input type="checkbox"/>
Pochodzi z korpusu IMPACT:	<input type="checkbox"/>
Tekst chroniony prawami autorskimi:	<input type="checkbox"/>

Znakowanie stylistyczno-genologiczne

Typ mowy: * niewierszowana ▾

Rodzaj: * teksty naukowo-dydaktyczne lub informacyjno-poradnikowe ▾

Gatunek:

- pieśni
- fraszki i epigramaty
- epitafia
- satyry
- sielanki
- kazania
- pisma polityczne
- polemiki religijne
- mowy okolicznościowe
- traktaty
- dialogi
- pamiętniki
- kroniki
- relacje
- opisy podróży
- herbarze
- akta sejmikowe
- wilkierze
- księgi sądowe
- inwentarze
- rejestry
- diariusze sejmowe
- rozmówki do nauki języka
- podręcznik
- przysłowia
- kalendarze
- przewodniki
- żywoty świętych
- poematy epickie
- przypowieści, specula (zwierciadła)
- modlitwy
- panegiryk
- lamentsy
- emblematy
- przywileje
- konstytucje sejmowe
- bajki
- księgi liturgiczne

Tematyka:

- alchemia
- anatomia
- architektura
- astrologia
- astronomia
- biologia
- botanika
- budownictwo
- chemia
- egzotyka
- ekonomia
- filozofia
- fizyka
- geografia
- gospodarstwo
- gramatyka
- górnictwo
- historia
- hutnictwo
- języki
- kulinaria
- matematyka
- medycyna
- mitologia
- miłość
- muzyka
- myślistwo
- obyczajowość
- poetyka
- polityka
- prawo
- religia
- retoryka
- wojskowość
- zielarstwo
- zoologia
- żeglarstwo

Poetyka żartu:



Znakowanie strukturalne i tekstowe

- Oznaczenie istotnych elementów struktury tekstu, np. elementów strony tytułowej, fragmentów obcojęzycznych itp.
- Znakowanie Korby odbywa się w szablonie dokumentu Word umożliwiającym wprowadzanie znaczników za pomocą skrótów klawiaturowych
- Na potrzeby projektu został stworzony program do konwersji plików Worda na format XML zgodny ze standardem TEI

Znaczniki strukturalne

[POCZĄTEK DOKUMENTU] - [KONIEC DOKUMENTU]
[POCZĄTEK STRONY] - [KONIEC STRONY]
[POCZĄTEK POZIOMU 1] - [KONIEC POZIOMU 1]
[POCZĄTEK STRONY TYTUŁOWEJ] - [KONIEC STRONY TYTUŁOWEJ]
[POCZĄTEK STRONY POTYTUŁOWEJ] - [KONIEC STRONY POTYTUŁOWEJ]
[POCZĄTEK MOTTA] - [KONIEC MOTTA]
[POCZĄTEK PRE-TEKSTU] - [KONIEC PRE-TEKSTU]
[POCZĄTEK POST-TEKSTU] - [KONIEC POST-TEKSTU]
[POCZĄTEK LISTY] - [KONIEC LISTY]
[POCZĄTEK INDEKSU] - [KONIEC INDEKSU]
[POCZĄTEK SPISU TREŚCI] - [KONIEC SPISU TREŚCI]
[POCZĄTEK NOTKI] - [KONIEC NOTKI]
[ILUSTRACJA]
[WZÓR MATEMATYCZNY]
[ZAPIS NUTOWY]
[INNA PRZERWA W TEKŚCIE]
[FRAGMENT NIECZYTELNY]
[TEKST W JĘZYKU OBCYM]
[ALFABET NIEŁACIŃSKI]

Znaczniki tekstowe

Tytuł publikacji	Podtytuł części
Powtórzony tytuł	<i>Streszczenie części</i>
Podtytuł publikacji	Podpis pod rysunkiem
<i>Autor</i>	Oznaczenie składki
Miejsce wydania	<i>Numer strony</i>
Drukarnia	<i>Kustosz</i>
<i>Rok wydania</i>	<i>Żywa pagina</i>
Inny element strony tytułowej	<i>Numer rozdziału</i>
<i>Tłumacz</i>	<i>Numer wersetu</i>
<i>Jezyk oryginału</i>	<i>Element dodany w Biblii Gdańskiej</i>
Numer wydania	Tekst odredakcyjny
Tytuł części	Tekst współczesny

Oznaczenia języków obcych

Arabski	Niemiecki
Czeski	Poludniowosłowiański
Francuski	Wschodniosłowiański
Grecki	Skandynawski
Hebrajski	Turecko-tatarski
Hiszpański	Węgierski
Litewski	Włoski
Łacina	Inny

Komentarze i poprawki skryptorów

- Informacja o liczbie nieczytelnych znaków
 - COMED<...>
- Informacja o możliwości błędnego odczytania słowa
 - gospodarz <?>
- Informacja o zaskakującym zapisie
 - gorzski <!>
- Rozwinięcie skrótów oznaczonych tyldą
 - atramēt <atrament>
- Poprawka oczywistej literówki
 - wszzędze <wszędzie>

Próbka tekstu – Word

[POCZĄTEK·DOKUMENTU]¶

[POCZĄTEK·STRONY·TYTUŁOWEJ]¶

PIECHOTNE·CWICZENIE¶

Albo· WOIENNOSC· PIESZA· Ktorą· Láćinnicy·
Pedestrem· Militiam· názywáią· Wodzom,
 Pułkownikom,· y· wszelkiew· Woienney·
Stárszynie·· Lubo· ktoszkolwiek· iest· Woienney·
spráwy··Miłóśnikiem.¶

Do·wiadomości·podána·¶

Przez·BLAZEIA·LIPOWSKIEGO¶

Nakładem·Ierzego·Forstera··I.K.M·Bibli·opoli·<?>¶

W·KRAKOWIE¶

W·Drukárni·v·Wdowy·Lukaszá·Kupiszá··I.K.M.¶

Typogr··Roku·Panskiego·1660.¶

[KONIEC·STRONY·TYTUŁOWEJ]¶

Próbka tekstu – XML

```

- <teiCorpus>
  <xi:include href="KORBA_header.xml"/>
- <TEI>
  <xi:include href="header.xml"/>
- <text>
  - <front>
    <pb/>
  - <titlePage xml:id="txt_1-titlePage">
    - <docTitle xml:id="txt_1.1-docTitle">
      <titlePart type="main" xml:id="txt_1.1.1-titlePart">PIECHOTNE CWICZENIE</titlePart>
      - <titlePart type="sub" xml:id="txt_1.1.2-titlePart">
        Albo WOIENNOSC PIESZA Ktorą Lácinnicy
        <foreign xml:lang="la" xml:id="txt_1.1.2.1-foreign">Pedestrem Militiam</foreign>
        nazywáią. Wodzom, Pułkownikom, y wszelkiej Woienney Stárszynie. Lubo ktoszkolwiek iest Woienney sprawy Miłóśnikiem.
      </titlePart>
    </docTitle>
    <note xml:id="txt_1.2-note">Do wiadomości podána.</note>
    <docAuthor xml:id="txt_1.3-docAuthor">Przez BLAZELA LIPOWSKIEGO</docAuthor>
  - <note xml:id="txt_1.4-note">
    Nakładem Ierzego Forstera, I.K.M. Bibli
    <uncertain reason="unreadable" xml:id="txt_1.4.1-uncertain">opoli.</uncertain>
  </note>
  - <docImprint xml:id="txt_1.5-docImprint">
    <pubPlace xml:id="txt_1.5.1-pubPlace">W KRAKOWIE</pubPlace>
    <publisher xml:id="txt_1.5.2-publisher">W Drukárni w Wdowy Lukaszá Kupiszá, I.K.M. Typogr.</publisher>
    <docDate xml:id="txt_1.5.3-docDate"> Roku Panskiego 1660.</docDate>
  </docImprint>
</titlePage>

```

Znakowanie morfosyntaktyczne

- Przypisanie każdemu segmentowi informacji gramatycznej:
 - o jego przynależności do określonej klasy gramatycznej
 - o wartościach kategorii gramatycznych
- Umożliwia zaawansowane przeszukiwanie korpusu, np. wyszukanie wszystkich form fleksyjnych leksemu
- Znakowanie Korby będzie się odbywać na wzór znakowania NKJP
- Tagset (zbiór znaczników) będzie oparty na tagsecie NKJP, który zostanie dostosowany do potrzeb opisu polszczyzny barokowej
- Znakowanie będzie się odbywać na tekstach transkrybowanych, a nie transliterowanych

Konwersja do wersji transkrybowanej

transliteracja	transkrypcja
<p>B. Nie z^áda mi żaden te^y rzeczy, ktoraby b^{ár}zi^ey owemu niżeli mnie należeć nie mi^áła. A t^ák śmierć ^y Kupidyn poz^árszy się i^áko bestie, posnęli w Koś^ćiele Bacchusowym, g^dzie niewyszumiawszy z przepi^ćia śmierć Cupidin^{ow}, ^á Cupido śmier^ći należący przypasawszy s^áy^dak do boku szły, swych należących odpr^áwow^áć powinnoś^ći.</p>	<p>B. Nie za^da mi żaden te^j rzeczy, któ^raby b^{ar}zy^ej owemu niżeli mnie należeć nie mi^ała. A ta^k śmierć ⁱ Kupidyn poz^arszy się j^ako beztie, posnęli w Koś^ciele Bacchusowym, g^dzie niewyszumiawszy z przepi^cia śmierć Cupidyn^{ów}, ^a Cupido śmier^ci należący przypasawszy sa^jdak do boku szły, swych należących odpra^wowa^ać powinnoś^ci.</p>

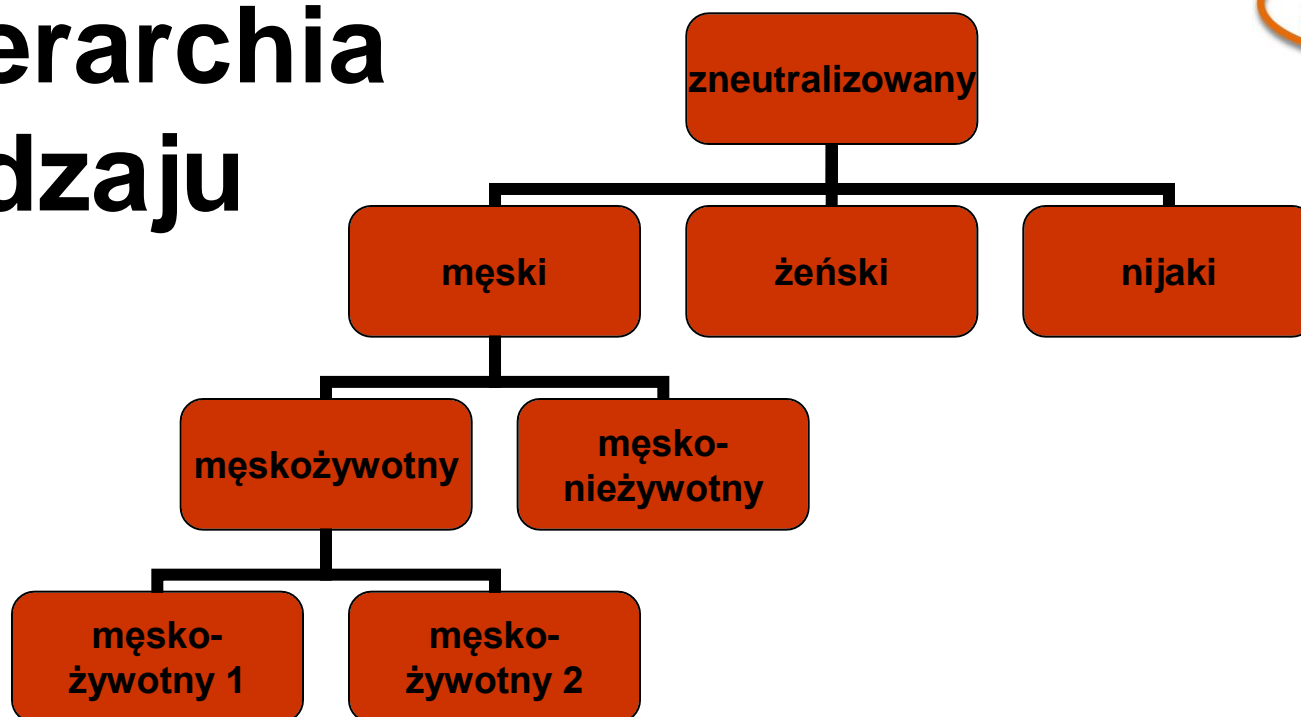
Narzędzia wspomagające znakowanie morfosyntaktyczne

- Analizator morfologiczny – przypisuje formie gramatycznej wszystkie możliwe interpretacje
- Tager – ujednoznacza interpretację gramatyczną
- Słownik gramatyczny – zawiera paradygmaty leksemów danego języka, stanowi podstawę działania analizatora morfologicznego
- Podstawą Morfeusza – analizatora morfologicznego przeznaczonego do analizy tekstów współczesnych (używanego przy znakowaniu NKJP) – jest *Słownik gramatyczny języka polskiego* (SGJP – <http://sgjp.pl>)
- Podstawą Morfeusza XVII-wiecznego będzie „postarzony” SGJP

„Postarzenie” SGJP

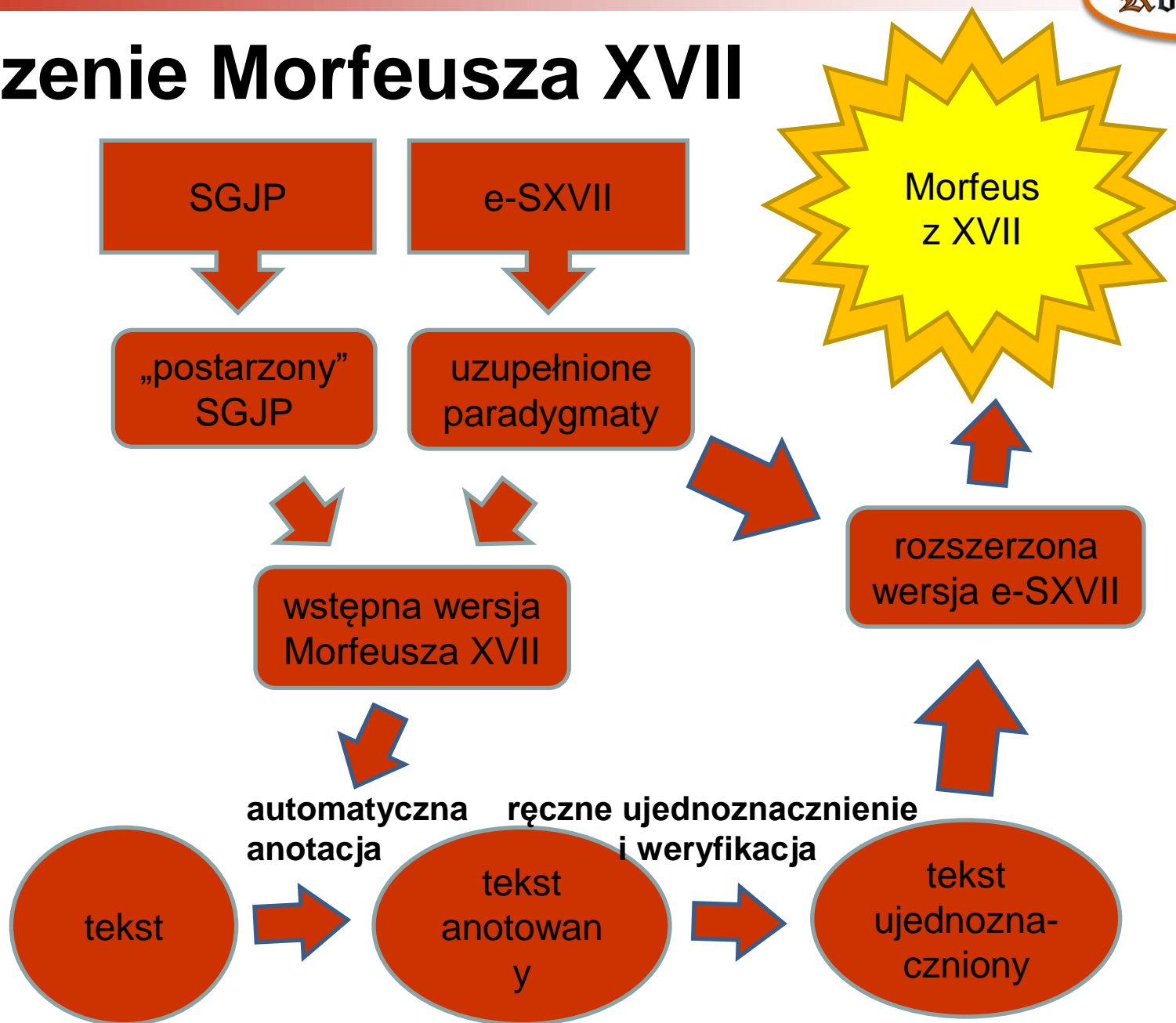
- Uzupełnienie SGJP o paradygmaty leksemów „wyciągniętych” z e-SXVII (niepełne paradygmaty zostaną automatycznie uzupełnione)
- Modyfikacja paradygmatów SGJP
 - Dodanie form historycznych do paradygmatów współczesnych, np. imiesłów przysłówkowy uprzedni zostanie dodany do paradygmatów czasowników niedokonanych (*widziawszy*)
 - Dodanie dodatkowych wartości do niektórych kategorii gramatycznych, np. liczby podwójnej (*żabie*)
 - Dodanie nowych, regularnie tworzonych derywatów, np. imiesłówów przymiotnikowych w stopniu wyższym (*bolątszy*)
 - Modyfikacja kategorii rodzaju

Hierarchia rodzaju



forma	wartości gramatyczne	lemat	rodzaj
abrysom	Dat. pl	ABRYS (?) ABRYSA (?) ABRYSO (?)	zneutralizowany
abrysu	Gen. sg	ABRYS	męski
abrys	Acc. sg	ABRYS	męskonieżywotny
abrysę	Acc. sg	ABRYSA	żeński

Tworzenie Morfeusza XVII





Dziękujemy za uwagę