

Format dokumentów
w projekcie elektronicznego korpusu
tekstów polskich z XVII i XVIII w.

Maciej Ogrodniczuk | Institute of Computer Science
Michał Lenart | Polish Academy of Sciences



INSTYTUT PODSTAW INFORMATYKI
POLSKIEJ AKADEMII NAUK
ul. Jana Kazimierza 5, 01-248 Warszawa

Instytut Języka Polskiego PAN

21 listopada 2013

Podobnie jak w NKJP, teksty w korpusie KORBA wraz z ich strukturą (podział na rozdziały, akapity, strony, przypisy, notki marginesowe itp.) znajdują się w plikach `text_structure.xml`.

```
<?xml version="1.0" encoding="UTF-8"?>
<teiCorpus xmlns="...">
  <xi:include href="KORBA_header.xml"/>
  <TEI>
    <xi:include href="header.xml"/>
    <text>
      <front> <!-- strona tytułowa i potytułowa -->
      <body> <!-- tekst właściwy - pozostałe strony -->
    </text>
  </TEI>
</teiCorpus>
```

```
<front>
  <titlePage>
    ...
  </titlePage>
  <div type="postTitlePage">
    ...
  </div>
</front>
```

Użycie poszczególnych elementów TEI z przykładami:

<http://www.tei-c.org/release/doc/tei-p5-doc/en/html>.

```
<titlePage>
  <docTitle>
    <titlePart type="main">The DUNCIAD,
                                VARIOURVM.</titlePart>
    <titlePart type="sub">WITH THE PROLEGOMENA
                                of SCRIBLERUS.</titlePart>
    <docAuthor>Gal Anonim</docAuthor>
  </docTitle>
  ...
</titlePage>
```

```
<titlePage>
  ...
  <docImprint>
    <publisher>Jan Karcan</publisher>
    <pubPlace>Wilno</pubPlace>
    <docDate>1611</docDate>
  </docImprint>
  <docEdition>Wydanie pierwsze</docEdition>
  ...
</titlePage>
```

Poszczególne części tekstu są reprezentowane przez elementy `<div>`; typ części znajduje się w atrybucie `@type`.

```
<body>
  <div type="volume">
    <p>...</p>
    <div type="part_level1">
      <!-- Rozdział 1 -->
      <p>...</p>
    </div>
    <div type="part_level1">
      <!-- Rozdział 2 -->
      ...
    </div>
  </div>
</body>
```

```
<head type="title">O nieporządnej chćiwości  
iedzenia,</head>
```

```
<head type="subtitle">to iest, o apetićie  
przećiwnym ludzkiemu  
przyrodzeniu.</head>
```

```
<head type="summary">POETA idąc zá zdánien dawnych  
Pogan, trzy poczátki świátá,  
y rzeczy będących ná nim,  
bydz rozumiał: Bogá, Lepsze  
przyrodzenie, y Chaos.</head>
```

Wartość atrybutu	type	Znacznik strukturalny
volume		[POCZĄTEK/KONIEC TOMU]
part_level	N	[POCZĄTEK/KONIEC POZIOMU N]
motto		[POCZĄTEK/KONIEC MOTTA]
preface		[POCZĄTEK/KONIEC PRE-TEKSTU]
epilogue		[POCZĄTEK/KONIEC POST-TEKSTU]
index		gdy zawiera indeks
contents		gdy zawiera spis treści
Wartość atrybutu	n	wczytywane z tekstu znajdującego się wewnątrz znacznika [POCZĄTEK/KONIEC CZĘŚCI NUMEROWANEJ]

Aktualnie w NKJP nie ma możliwości sprawdzenia, na której stronie tekstu źródłowego znajduje się znaleziony fragment.

Aby to umożliwić, do dokumentów `text_structure.xml` został włączony znacznik `<pb>` (ang. *page break*) oznaczający początek nowej strony.

Atrybut `n` zawiera numer strony lub inne oznaczenie identyfikujące.

Początek dokumentu zawiera oznaczenie strony pierwszej:

```
<front>
  <!-- początek strony tytułowej -->
  <pb n="1"/>
  <titlePage>
    ...
  </titlePage>
</front>
```

```
<body>
  <div type="volume">
    ...
    <p>Womit nic inszego nie iest/
      iedno wyrzucanie przez vstá tego/
      coby iedno w nim było. Cknienie zaś
      iest początek niepewnego
      <pb n="2"/>
      poruszenia mocy wyganiáiącey/
      przez którą vsięluie przez vstá wygnác/
      to co iest żołądkowi przykro/
      ábo nieznośno.</p>
    <p>...</p>
    ...
  </body>
```

Uwaga:

- Oznaczenie kustoszy jest przepisywane dla celów weryfikacyjnych, ale nie jest reprezentowane w formacie wynikowym.
- Jeśli słowo znajduje się na dwóch stronach jednocześnie (czyli zawiera przeniesienie), wówczas znacznik końca strony zostanie umieszczony za tym słowem.
- Każda strona zaczyna się znacznikiem <pb> — nawet strona tytułowa i potytułowa.

Śródczęściowy podział na strony

```
<body>
  ...
  <div type="volume">
    ...
    <div type="part">
      <head>Rozdział 1</head>
      ...
    </div>

    <!-- początek strony bezpośrednio w <div>,
         czyli np. pomiędzy częściami -->
    <pb n="12"/>

    <div type="part">
      <head>Rozdział 2</head>
      ...
```

W notkach marginesowych zapisywane są różne fragmenty, które nie stanowią treści tekstu, ale które chcemy zachować.

```
<p>Wielkość zamyka w sobie/  
wielkie ábo nie wielkie iedzenie/  
y częste iedzenie. Co się tknie pierwszego.  
<note type="margin">Iák wiele ieść</note>  
Iák wiele mamy ieść/  
naznacza nam všmierzenie głodu/  
to iest/ gdy się zniešie przez  
wzięty pokarm przyrodzona chćiwość iadłá.</p>
```

Uwaga: narzędzia automatycznie przetwarzające teksty muszą szczególnie uważać na ten znacznik, ponieważ jego zawartość może „zaśmiecać” właściwy tekst (notka marginesowa i przypis są jakby drugim wymiarem „płaskiego”, jednowymiarowego tekstu).

Przypisy:

```
<p>Y rzekł Mojżesz przed Pánem/  
mowiąc: Oto/  
Synowie Izráelscy nie usłucháli mię/  
á jákoż mię usłucha Fáráo/  
á jam nie obrzezánnych  
<note type="annotation">wolney wymowy nie mam.</note>  
warg?</p>
```

Podpis pod rysunkiem:

```
<note type="label">a. To iest tępy koniec kliniká,  
ktory ma byđź zewnątrz w vlu.</note>
```

Żywa pagina:

```
<fw type="runningHead">0 Przypadkách Brzemiennych.</fw>
```

Znacznik strukturalny	Znacznik TEI
[ILUSTRACJA]	<code><figure/></code>
[WZÓR MATEMATYCZNY]	<code><formula/></code>
[ZAPIS NUTOWY]	<code><notatedMusic/></code>
[TEKST W JĘZYKU OBCYM]	<code><foreign xml:lang="kod"/></code>
[ALFABET NIEŁACIŃSKI]	<code><foreign xml:lang="x-nlws"/></code>
[INNA PRZERWA W TEKŚCIE]	<code></code>
Znak specjalny	<code><g/></code>

Język	Kod
arabski	ar
czeski	cs
francuski	fr
hebrajski	he
hiszpański	es
litewski	lt
łacina	la
niemiecki	de
południowo-słowiański	zls
ruski	x-ruski
skandynawski	x-skand
turecko-tatarski	x-turtat
węgierski	hu
włoski	it

Oznaczenia rozpoczynające się literą x pochodzą spoza oficjalnej listy kodów językowych ISO; to tagi językowe do użytku prywatnego, wykorzystywane tam, gdzie standard okazuje się niewystarczający.

Z wewnętrznej korespondencji:

Zdecydowaliśmy się połączyć kilka języków w jedną grupę, bo w XVII w. niektóre języki albo jeszcze nie istniały, albo byłyby trudne do odróżnienia.

Turecko-tatarski obejmuje turecki i tatarski, skandynawski — teoretycznie wszystkie języki skandynawskie, ale praktycznie szwedzki i duński. Zdecydowaliśmy się też nie używać określenia „ukraiński”.

Opisany format zapisu jest oparty na NKJP, ale wprowadza wiele nowych znaczników. Narzędzia oparte na NKJP (np. Poliqarp, Pantera) obsługują tylko uproszczoną wersję formatu.

Wymagane jest dostosowanie wyżej wymienionych narzędzi lub stworzenie zubożonej wersji formatu.

W ramach projektu powstaje serwis internetowy dla skryptorów i redaktorów.

Umożliwia on między innymi:

- konwersję plików .doc do formatu TEI za pomocą formularza na stronie,
- wykrywanie błędów w oznaczaniu tekstu (brakujące znaczniki, źle oznaczone języki obce itp.),
- podgląd wyników konwersji w formie czytelnej dla użytkownika (ułatwia znajdowanie błędów),
- zarządzanie procesem dostarczania tekstów — wgranie wstępnej wersji przez skryptora, przesyłanie poprawek, redakcja (jeszcze nie zaimplementowane).

KORBA - wczytywacz tekstów

[Strona główna](#) Witaj **mLenart** | [Administracja](#) | [Zmień hasło](#) | [Logout](#)

Plik .doc: No file chosen

Lista plików

ŻędzKom-poprawione-mLenart.doc	wynik konwersji	czytaj
próbka_KGO_1678-poprawione-mLenart.doc	wynik konwersji	czytaj
PetrSInst-mLenart.doc	wynik konwersji	czytaj

KORBA - wczytywacz tekstów

[Strona główna](#) [Witaj mlenart](#) | [Administracja](#) | [Zmień hasło](#) | [Logout](#)

- Błąd przy wczytywaniu pliku: Znacznik zamknięcia nieotwartej strony.

Strona 36

ze WLADYSŁAWA IV spoiony dnia 13 Września rękoma
Arcybiskupiami w Krakowie zwielkim światá wszytkiego
weselem/ Państw/ Monárcy życzliwością Koronowana.
záćmił/F2[KONIEC STRONY]

[KONIEC STRONY]

46

ORATOR

nie ostatnia/ największe po Krolewsku do Niebá iść. Nászá
Cecilia Renatá y po Krolewsku się vrodziłá/ y po Krolewsku
żyłá/ y po Krolewsku vmárlá/ y po Krolewsku do N

Plik .doc: No file chosen

KORBA - wczytywacz tekstów

[Strona główna](#) Witaj **mLenart** | [Administracja](#) | [Zmień hasło](#) | [Logout](#)

- Dodano dokument PetrSInst-mLenart.doc, ale prawdopodobnie były błędy
- Prawdopodobnie źle zaznaczony język obcy

R.Vnguenti sandalini ź 1.
Pul Pimpinellae
Scordii
Mirrhae an -) 1
Succi Citri žiiŷ
Misce et f. linimen

- Prawdopodobnie źle zaznaczony język obcy

R.Vnguenti sandalini ź 1.
Pul Pimpinellae
Scordii
Mirrhae an -) 1
Succi Citri žiiŷ
Misce et f. linimentum pro corde.

Plik .doc: No file chosen

Strona tytułowa	
Tytuł	INSTRVCTIA Abo NAVKA/ IAK SIĘ SPRAWOWAC CZASV moru.
Podtytuł	W ktorey się zamyka: 1. Ochroná: Iáko się vchrániác morowego powietrza. 2. Leczenie wsztykich niemal przypadków w nim/ gdzieby kogo opanowało. Dla prostych nápisána/ krom discursow:
Miejsce wydania	W KRAKOWIE,
Drukarnia	W Drukárni Mikołáiá Lobá/
Rok wydania	Roku Páńskiego/ 1613.
Autor	Przez DOC. SEBASTYANA PETRYCEGO MEDYKA.

postTitlePage
motto
Trway w sławie/ poty/ nie<>obyta Wieży/ Dokąd zołw wsztykiego świata nie obieży: Abo dokąd mrowka morzá nie wypiie/ Lub dokąd vmárly znouw nieożyie.

preface
Szláchetnym y zacnym Pánom/ ICH MCIOM: P. KASPROWI CVTETEROWI B. IOAHIMOWI

materiał na wierzch ze wnętrza być wyrzuconą. Tak poznawszy powietrze teraz o tym będziem mówić/ iak się go vchronić przez dietę/ y pewne lekárstwa. A potym iak go leczyć ma/ kto się go vchronić nie mogł.

Ochroná powietrza przez porzádne życie.

Przestogá generalna 1.

Nim o rzeczách do ochrony zarázy szczegulnie mówić pocznę/ napotrzebniejszą powszechnią przestrogę naprzod polożę. Do Páná Bogá się vćiecz z pokutą y próżbą.

Pan Bog przyczyna moru

O powietrzu.

Bo on sam taką plagę na ludzi dopuszcza: bez iego woli nie stáie się nic. Zadna ochroná nie pomoże/ gdy iego łaski y straży nad námi nie bedzie. Ptacy/ zá ie^o^ rzádem y staraniem swoje żywność máia/ á cożby ludzie nie mieli mieć: Ten na ludzie przepuszcza vporne choroby/ y přetkie zarazy: ten od nich broni. A to ieszcze dziwney táiemnicy Páńskiej przypisác potrzebá/ iz czasem choc się kto nawięcey nátrąca na nie/ minie y nie imie się go: a podczas tego/ co nie dáie przyczyny/ vchwyći. Co Doktorowie niektorzy doświadczyli naa sobie/ ktorzy choc chodźili do zaráżonych/ niezárázili siebie samych/ áni od siebie drugich. Rychley na tego padło co się nawięcey chronił.

Przestoga generalna 2.

Druga powszechnia przestogá niech będzie. Ze złego powietrza/ na dobre vbiegác. Co y nieme zwierzęta czynią: pod taki czas vćiekaią z iam swoich/ á inszych sobie szukaią zdrowych. A nie