

## Разметка корпуса Азбуковника 1596 г. на основе рекомендаций TEI

### Digitalising of the Russian Manuscript Dictionary of 1596 Using the TEI Encoding Scheme

Коваленко К.И., магистр, мнс

kira.kovalenko@gmail.com

Институт лингвистических исследований

Российской академии наук,

Санкт-Петербург, Россия

Ключевые слова: историческая лексикография, азбуковники, разметка текстов, TEI

The paper deals with the problems of Russian manuscript dictionary digitalising. The dictionary has a complicated entry structure, which combines different types of information: word origin, definitions, quotations from literary texts, references to other words, etc. It is supposed to use the TEI encoding scheme, as it provide a wide variety of tags for manuscripts and dictionaries encoding. The digitalising would provide more effective access to the information and accelerate linguistic researches.

Создание корпусов текстов — одно из наиболее перспективных направлений современной лингвистики. В настоящее время в нашей стране активно развивается проект Национального корпуса русского языка, одно из направлений которого предполагает расширение корпуса за счет включения текстов, отражающих историю русского языка, — от берестяных грамот и летописей до текстов XVIII — начала XIX вв. В частности, предполагается создать Корпус старорусских текстов XV-XVII вв., в рамках которого ожидается пополнение корпуса непосредственно за счет текстов в электронном виде, а также разработка компьютерной морфологии и словаря старорусских текстов [<http://www.corpling-ran.ru/n2.html>]. Можно предположить, что корпус средне-русских текстов стал бы более репрезентативным, если бы в него были вклю-

чены тексты азбуковников XVI-XVII вв. — предшественников современных словарей.

Азбуковник как лексикографический жанр сложился к середине XVI в. на основе более ранних словарных сводов — ономастиконов (гlossариев имен библейских персонажей и наименований библейских мест), приточников (перечней символов Псалтири), glossариев к «Лестнице» Иоанна Синайского и словарей-разговорников [История русской лексикографии 2001: 40-49]. В дальнейшем азбуковники активно пополнялись новыми статьями, источником которых служили glossы на полях рукописей, пояснения непонятных читателю реалий в текстах литературных произведений. Поэтапный характер формирования азбуковников, разнородные сведения, представленные в нем, являются причиной того, что статья азбуковника включает в себя достаточно разнородную информацию и далеко не всегда имеет четко выстроенную линейную структуру.

Азбуковник, выбранный для разметки, является типичным представителем своего жанра. Он был создан в Новгороде в 1596 году в монастыре Антония Римлянина и сохранился до настоящего времени в единственном списке начала XVII в. (РНБ, собр. Погодина, № 1642). В нем содержится более 6000 статей, которые организованы по первой букве и по следующей гласной. В отличие от своих предшественников, в нем представлена не только лексика греческого, латинского, старославянского и тюркского происхождения, но имеется значительный пласт лексических единиц из польского языка.

Структура статьи азбуковника имеет достаточно сложную организацию. Так, в заголовочной части статьи может быть слово, словосочетание или целая фраза (например: *Кеси ерѳтисо<sup>h</sup> апаоуѳонъ* или *Что е<sup>c</sup> пать чювьствъ дшеѳвны<sup>x</sup> и пать тѳлѳсныхъ*). Зона толкования может содержать объяснение сразу нескольких слов (*Житіе и жизнь (т) житіе нарицае<sup>m</sup> еже кто какова ѿца и коегѳ града и коеѳ вѳры . жѳзнь же се е<sup>c</sup> каковыми дѳлы бгѳ у оугодѳ, і койми дарованіи ѿ гѳ прославле<sup>h</sup>, и како теченіе подвига сконча*), указание на паронимы (*Граногра<sup>o</sup> (т) гранѳи писе<sup>u</sup> землемѳрны<sup>x</sup>. іно бо е<sup>c</sup> граногра<sup>o</sup>. і іно хроногра<sup>o</sup>*), а также дополнительную информацию, относящуюся к толкованию данного слова опосредованно (*Всеплѳдіе (т) вѳ ветхо<sup>m</sup> законѳ прѳведе<sup>o</sup>*

*о́вча. и́ли во́ль на жрѣ́тву въ ѣ́ви. е́гда закла́вше всесо́жгу<sup>m</sup>, то́ всепло́дие и́менуе<sup>m</sup>сѧ. а́ е́гда закла́вше, нѣ́кіѧ ча́сти бгѣ́у ѡ́дѣла<sup>m</sup> нѣ́кіѧ же і́ереѡ<sup>m</sup> і́ наро́до<sup>m</sup> на снѣ́деніе, то нари́че<sup>m</sup>сѧ жѣ́ртва).* Как правило, над заготовочным словом указывается предполагаемый язык-источник, однако данная помета может присутствовать и над любым другим словом в любой части статьи (*Салафѡѧ {ж}*, *и́рѣна {г}*, *ѣже е<sup>c</sup> мирнаѧ*). В некоторых случаях на полях указываются литературный источник (источники), из которых данная статья попала в азбуковник. Но иногда ссылка на источник оказывается без пометы, и остается только догадываться, к какой же статье азбуковника она относится.

Разметку корпуса словаря предполагается осуществить на основе рекомендаций TEI — Text Encoding Initiative (версия P5). На базе рекомендаций TEI уже были реализованы и продолжают функционировать более 150 проектов [<http://www.tei-c.org/Activities/Projects>]. Среди отечественных проектов можно отметить информационно-поисковую систему "Русская литература XVIII века" [<http://antology-xviii.spb.ru>], также выдвинут ряд предложений по использованию TEI в работах [Вадряев 2005; Вотинцев 2006; Бабалык, Варфоломеев, Пигин 2010; Захаров, Митрофанова, Михайлова 2011].

Все более широкое распространение стандартов TEI обусловлено тем, что рекомендации по разметке рассчитаны на разные типы текстов: прозаические и стихотворные произведения, словари, транскрипцию разговорной речи. Отдельное внимание уделяется электронному представлению текстов рукописных источников. Необходимо также отметить, что рекомендуемая разметка рассчитана в том числе и на словари со сложной структурной организацией, что дает возможность пользоваться исключительно предложенными тегами и атрибутами, не усложняя код разметки дополнительными нововведениями.

В качестве примера использования предлагаемых TEI тегов возьмем статью азбуковника *Кидарь {ж} {ѡлѡ<sup>m</sup> рѡі} (т) тма. и́ли кло́букъ а́рхіе́реѡскіѡ. и́же и́ митра глѣ́тсѧ. і́ па́ки кидѧ<sup>p</sup> нари́цае<sup>m</sup> и́ до́мы татарскіѧ. о́вы бо и́хъ в землѣ и́ско́паны , и́ того радѣ́ темныѡ; о́вы н а колесни́ца<sup>x</sup>, по<sup>o</sup>стмѡ о́граждѣны; 'ѣсть же и́ гра<sup>o</sup> нари́цае<sup>m</sup> кидарь о́ немже пи́шетъ , вселѣ́хсѧ с с елы кидарьскіѡми.* Разметка данной статьи, ориентированная на структурную организацию, будет выглядеть следующим образом:

<entryFree>

<form><w>Кидáръ<note><lang>ж</lang><bibl><title>ѡлѡм рѣ</title></bibl></note></w></form>

<sense> (т) <def n="1">тма.</def> и́ли <def n="2"> клобукъ а̀рхієрѣискіи. ѡ́же и́ <mentioned>митра</mentioned> глѣтса.</def> Ї па́ки кидар нарица́ет и́ <def n="3">до́мы татарс́кіа. о́вы бо ѡ́хъ в земли́ и́скопаны, и́ того рад“ тѣмныи; о́вы на колесни́цах, полстмй о́граждѣны;</def> <def n="4">’ѣсть же и́ град нарица́емъ кидáръ о́ немѣже пи́шетъ, <cit> <quote>вселѡхса с селы кидарьскіи</quote> </cit></def></sense>

</entryFree>

Кроме разметки текста в соответствии со структурой статьи, предполагается ввести дополнительные сведения: по возможности определить иноязычное слово в языке-источнике, лексикографический или литературный источник и контекст, послуживший материалом для создания той или иной статьи, тематическая разбивка лексики по группам. Это должно значительно облегчить лингвистическое исследование корпуса азбуковника, в частности, поможет точнее определить, из каких языков лексика той или иной тематической группы наиболее активно проникала в русский язык, а также выявить фонетические и графические приемы передачи иноязычных слов.

### Литература

Бабалык, Варфоломеев, Пигин 2010 — Бабалык М.Г., Варфоломеев А.Г., Пигин А.В. Использование формата TEI для публикации и анализа списков произведений вопросо-ответного жанра // Информационные технологии и письменное наследие: материалы междунар. науч. конф. (Уфа, 28-31 октября 2010 г.) / отв. ред. В.А. Баранов. Уфа; Ижевск: Вагант, 2010. С. 17-20.

Вадяев С.Е. Лингвистические принципы построения и использования корпуса текстов для исследования официально-делового стиля современного немецкого языка (на материале электронного корпуса «DER»). Автореф. дис. на соиск. учен. степ. к.филол.н. Н. Новгород, 2005.

Вотинцев П.А. Использование формата TEI для обмена данными с полнотекстовой информационно-поисковой системой «Манускрипт» // Современ-

ные информационные технологии и письменное наследие: от древних рукописей к электронным текстам: материалы междунар. науч. конф. (Ижевск, 13-17 июля 2006 г.) / отв. ред. В.А.Баранов. Ижевск: Изд-во ИжГТУ, 2006. С. 30-31.

Захаров, Митрофанова, Михайлова 2011 — Захаров В.П., Митрофанова О.А., Михайлова В.Д. Разметка словарей в соответствии со стандартом TEI // Информационные технологии в лексикографии. СПб.: СПбГУ, Филологический факультет, 2011. С. 61-69.

История русской лексикографии 2001 — История русской лексикографии / отв. ред. Ф.П. Сороколетов. СПб.: Наука, 2001.