



Писменото наследство и информационните технологии

El'Manuscript – 2014

МАТЕРИАЛИ ОТ V МЕЖДУНАРОДНА НАУЧНА КОНФЕРЕНЦИЯ
ВАРНА, 15–20 СЕПТЕМВРИ 2014 Г.

София • Ижевск
2014

Специализированный лингвистический корпус социокультурной специфики регионального варианта русской речи¹

Л. В. Рычкова

Русский язык, региональный вариант русской речи, социокультурная специфика речи, лингвистический корпус, корпус СМИ

A Specialized Linguistic Corpus Created for Research on the Sociocultural Characteristics of a Regional Variant of Russian Speech

Ludmila Rychkova

This paper describes peculiarities of the development of a specialized linguistic corpus comprising Russian-language texts from the mass media printed in the Minsk Region of Belarus, the goal of which is research on the sociocultural characteristics of one of the regional varieties of Russian speech.

Переход от исследования собственно лингвистических реалий к использованию корпусных данных для выявления комплекса социокультурных факторов, влияющих на варьирование языка и обуславливающих принципиальную изменчивость и вариативность функционирования его единиц, стал возможным благодаря формированию новой, экспериментально-доказательной, парадигмы лингвистики, становлению которой, без сомнения, способствует лингвистика корпусная, в рамках которой не только создаются такие репрезентативные массивы корпусных данных, как лингвистические корпуса текстов, но и разрабатываются методики и направления междисциплинарных исследований, основанные на обобщении корпусных данных.

Процесс формирования лингвистических корпусов — это сложная и трудоемкая процедура, в которой важную роль играет грамотное определение этапов формирования корпуса, обеспечение аутентичности отобранного в качестве базы корпуса языкового материала, разработка системы метаразметки, в том числе — тематической. Учет всех этих факторов позволяет обеспечить информационно-лингвистическую релевантность корпуса, достоверность корпусных данных. Степень информативности корпуса зависит также от использованных в нем видов лингвистической разметки. Автоматический режим разметки, как правило, дает много «информационного шума», обусловленного принципиальной неоднозначно-

¹ Выполнено в рамках проекта, реализуемого при поддержке Белорусского республиканского фонда фундаментальных исследований (договор № Г13Р-050 от 16.04.2013 г.).

стью языковых объектов. В связи с этим первично размеченный корпус до ввода в эксплуатацию пройти этапы выверки и апробации.

В рамках международного белорусско-русского проекта создается историко-лингвистический корпус нового типа, представляющий собой оригинальный электронный языковой ресурс, специфика которого обусловлена, прежде всего, особенностями избранного в качестве основы для создания данного корпуса материала.

В частности, важной особенностью создаваемого корпуса, которая проявляется в процессе сбора и систематизации электронного контента, является его «решающий» характер: в корпусе сохраняются и специальным образом размечаются белорусскоязычные фрагменты и тексты, отражающие естественное в условиях близкородственного белорусско-русского двуязычия сознательное переключение кодов, достаточно часто осуществляемое в текстах СМИ Гродненщины. Таким образом, такое сознательное переключение кодов не имеет ничего общего с так называемой «трясянкой», представляющей собой результат неосознанного психологического двуязычия смешанного типа). Все существующие до сих пор языковые корпуса создавались либо как параллельные, либо как сопоставимые образцы, создаваемый корпус представляет собой первый опыт «смешанного» корпуса, отражающего соответствующий языковой материал, что создаст основу для проведения межъязыковых исследований.

Второй важной особенностью создаваемого корпуса является его «региональный» характер, так как он строится на материале текстов СМИ Гродненщины полиэтнического региона Беларуси, непосредственно граничащего с Польшей и Литвой. Как отмечают российские участники проекта — непосредственные разработчики Национального корпуса русского языка (НКРЯ) [Национальный корпус русскоязычных СМИ Гродненщины положит начало формированию в составе НКРЯ нового модуля, представляющего региональные варианты русского стандарта] [Кустова, Савчук 2013: 352].

Поскольку «новая языковая ситуация в государствах постсоветского пространства приводит к усилению процессов дивергенции региональных вариантов русского языка и кодифицированного языка метрополии, к постепенному формированию национальных вариантов» [Кустова, Савчук 2013: 345], то создание такого корпуса является объективным отражением «генерального» характера НКРЯ, стремящегося отразить все функциональные варианты русского языка.

Выбор в качестве базы корпуса текстов СМИ не случаен. Именно СМИ наиболее «чувствительны» ко всем факторам социо- и лингвоэкологии и наиболее ярко отражают различные виды актуального дискурса, и поэтому позволяют получать данные, релевантные для выявления различного рода социокультурной специфики национально или регионально обусловленных идиомов. СМИ занимает сегодня главенствующее положение в системе этнической и межкультурной коммуникации, оказывая влияние на формирование стереотипов и нормативного речевого поведения и социокультурных конструктов. В текстах СМИ

корпус не только общеупотребительная лексика, но и лексика специальная, в частности, именно тексты СМИ отражают формирование языков для социокультурных целей. Динамичность языка СМИ позволяет выявлять семантический потенциал лексических единиц и специфику его реализации в различных типах обусловленную аксиологическими, гендерными и иными социокультурными факторами.

«Постративный» характер корпуса текстов СМИ Гродненщины, обусловленный вынужденным, по причине высокой степени трудоемкости создания корпусом ограничением объема исходного языкового материала, актуализировал необходимость отбора конкретных изданий, тексты которых составили базу для формирования «сырого» (неразмеченного) электронного языкового ресурса на первом этапе создания корпуса.

В частности, была осуществлена первичная характеристика СМИ с учетом следующих критериев: название газеты, сайт, место издания, официальность, неофициальность, год основания, периодичность, язык статей, адрес издательства, тематика. Как важные факторы в отборе СМИ для формирования лингвистического корпуса рассматривались следующие: наличие интернет-версий СМИ, их доступность и возможность конвертации в формат корпуса. Кроме того, был выполнен мониторинг интернет-версий русскоязычных СМИ Гродненщины, в результате чего был определен характер существующих архивов различных СМИ, их объем, пригодность для формальной и технической обработки для целей создания лингвистического корпуса.

Поскольку описываемый в данной работе корпус создается как экспериментально-показательная база для исследования социокультурной специфики русского языка Гродненщины, то особое внимание было уделено разработке системы метки, позволяющей должным образом структурировать исходный языковой материал в соответствии с релевантными целями исследования частными лингвистическими задачами.

В дальнейшем остановимся подробнее на двух параметрах метаразметки иллюстративного корпуса текстов СМИ Гродненщины — «Сфера функционирования» и «Тематика».

Параметр «Сфера функционирования» в массиве текстов корпуса представлен следующими значениями: «бытовая», «официально-деловая», «производственно-сервисная», «публицистика», «реклама», «учебно-научная», «художественная», «религиозно-богословская», «электронная коммуникация».

С первого взгляда, набор значений не слишком велик, но в любом случае, он отражает фактическое стилевое разнообразие текстов современных СМИ, выходящее за пределы собственно публицистики.

Одним из важных параметров системы метаразметки является тематика текстов СМИ. В газетном корпусе НКРЯ тематическая разметка не предусмотрена. Тем не менее, в основном корпусе НКРЯ все нехудожественные тексты размечены в соответствии с тематикой.

Определение тематики текста при осуществлении разметки вызывает определенные сложности, так как большие по объему тексты чаще всего, как правило, информационно неоднородны, то есть в них представлено несколько тем. Поскольку при осуществлении тематической разметки текстов нехудожественной литературы основного корпуса НКРЯ в случае неоднозначности дается список индексов, соответствующих всем предметным областям, нашедшим отражение в тексте, то аналогичным образом осуществлялась тематическая разметка и при создании иллюстративного корпуса русскоязычных СМИ Гродненщины. Для определения перечня адекватных фактическому языковому материалу тематических индексов в автоматическом режиме были получены списки тематических маркеров архивов газет — источников исходных для корпуса текстов. В ручном режиме была составлена таблица соответствий полученных таким образом тематических маркеров значениям параметра “тематика”, предложенной разработчиками НКРЯ для текстов нехудожественной литературы основного корпуса. Анализ полученных таким образом данных показал, что, с учетом допущения, этих индексов достаточно для тематической разметки текстов корпуса СМИ Гродненщины. В итоге было добавлено лишь одно значение параметра — “происшествие”.

Принципиальное соответствие значений параметров метаразметки иллюстративного корпуса текстов СМИ Гродненщины характеристикам основного корпуса письменных текстов НКРЯ — важный фактор, обеспечивающий “включение” этого корпуса в новый модуль НКРЯ — корпус региональной и зарубежной прессы.

Литература

- Кустова, Савчук 2013 — *Кустова Г. И., Савчук С. О.* Изучение лексико-семантической и социокультурной специфики русской речи на территории Республики Беларусь (на материале текстов СМИ) // Труды междунар. конф. “Корпусная лингвистика–2013”, Санкт-Петербург, 25–27 июня 2013 г. СПб.: С.-Петербургск. гос. ун-т Филол. фак-т, 2013. С. 344–352.
- Национальный — *Национальный корпус русского языка* [Электронный ресурс]. Режим доступа: <http://www.ruscorga.ru>, свободный (дата обращения: 17.02.2014).