

Национальный корпус русского языка как основа новаторских электронных учебников

Сибирцева Вера Григорьевна

к.ф.н., доцент кафедры прикладной лингвистики и межкультурной коммуникации факультета гуманитарных наук;

Национальный исследовательский университет – Высшая школа экономики
ул. Большая Печерская, д.25/12, г.Нижний Новгород, 603155, (831) 436-74-09
vsibirtseva@hse.ru

Хоменко Анна Юрьевна

магистрант 1 курса факультета Бизнес-информатики и прикладной математики

Национальный исследовательский университет – Высшая школа экономики
ул. Большая Печерская, д.25/12, г.Нижний Новгород, 603155, (831) 436-74-09
lili-th89@narod.ru

Баранова Юлия Николаевна

магистрант 1 курса факультета Бизнес-информатики и прикладной математики

Национальный исследовательский университет – Высшая школа экономики
ул. Большая Печерская, д.25/12, г.Нижний Новгород, 603155, (831) 436-74-09
ligros7@gmail.com

Аннотация

В статье идёт речь о разработках научно-учебной группы студентов и преподавателей Национального исследовательского университета – Высшей школы экономики «Корплинги (Нижний Новгород-Москва)». Данная работа связана с исследованиями в области компьютерной и корпусной лингвистики. Разработки нацелены в первую очередь на создание интерактивных ресурсов, основанных на материалах Национального корпуса русского языка, для обучения студентов русскому языку как иностранному. В статье отражены уже проведённые этапы работы научно-учебной группы: 1) решение задачи, связанной с русским глаголом и продуктивным префиксальным словообразованием, 2) решение задачи, связанной с адаптацией языкового материала Национального корпуса русского языка для базы данных электронного учебника «Русский язык как иностранный». В статье описан не только алгоритм решения поставленных задач и конечный результат исследований, но и трудности, с которыми столкнулись разработчики, а также пути их решения.

The article reports about the students and teachers research group of National Research University Higher School of Economics entitled "Corplingui (Nizhny Novgorod-Moscow)" development. This work is about the research in the field of computer and corpus linguistics. Development primarily focuses on the creation of interactive resources based on the materials of The Russian National Corpus. The interactive resources target to teach students Russian as a foreign language. The paper describes the stages of the work, which have been already carried out by the research group: 1) the solution of the problem related to the Russian verb and productive prefix word-formation; 2) the solution of the problem related to the adaptation of linguistic material of The Russian National Corpus for database of the electronic textbook "Russian as a foreign language". This paper describes not only the algorithm for solving the tasks and the final result of the research, but also the

difficulties, which the developers face, and their solutions.

Ключевые слова

русский язык как иностранный, Национальный корпус русского языка, компьютерная лингвистика, корпусная лингвистика, язык программирования Python

Russian as a foreign language, The Russian National Corpus, computational linguistics, corpus linguistics, Python programming language

Введение

В мировой практике корпусная педагогика (Corpus-Based Approach) завоевывает серьезные позиции, многочисленные научные работы, базирующиеся на лингвистических корпусах, посвящаются исследованию функционирования метафоры, особенностям лексического состава языка, анализу словообразовательных лингвистических возможностей. Корпусные технологии ощутимо влияют на современную лингводидактику, в которой обращение к параллельным корпусам становится нормой. Национальный корпус русского языка (далее – НКРЯ) является крупнейшим русскоязычным корпусом, но его дидактические возможности до настоящего времени не используются в обучении в полной мере, хотя технологично организованная опора на ресурс такого масштаба открывает перспективы создания учебных пособий совершенно нового типа: современных и минимально адаптированных.

В 2012 году в Национальном исследовательском университете – Высшей школе экономики (далее – НИУ ВШЭ) началась работа над подготовкой материала НКРЯ для создания электронных учебников по русскому языку как иностранному. Научно-учебная группа студентов и преподавателей НИУ ВШЭ «Корплинги (Нижний Новгород-Москва)» находится в настоящий момент в середине исследовательского и творческого пути, и данная статья посвящена рассмотрению этапов работы и анализу выявленных в процессе работы достоинств и недостатков как самого проекта, так и исследуемого материала.

Материалы Национального корпуса русского языка в грамматических упражнениях

На основе данных Национального корпуса русского языка и исследования моделей русского языка собран материал для электронного учебного пособия по теме «Русский глагол. Продуктивное префиксальное словообразование». В теоретической части пособия глагольные префиксы структурированы по значениям и снабжены примерами из художественной литературы и современной публицистики.

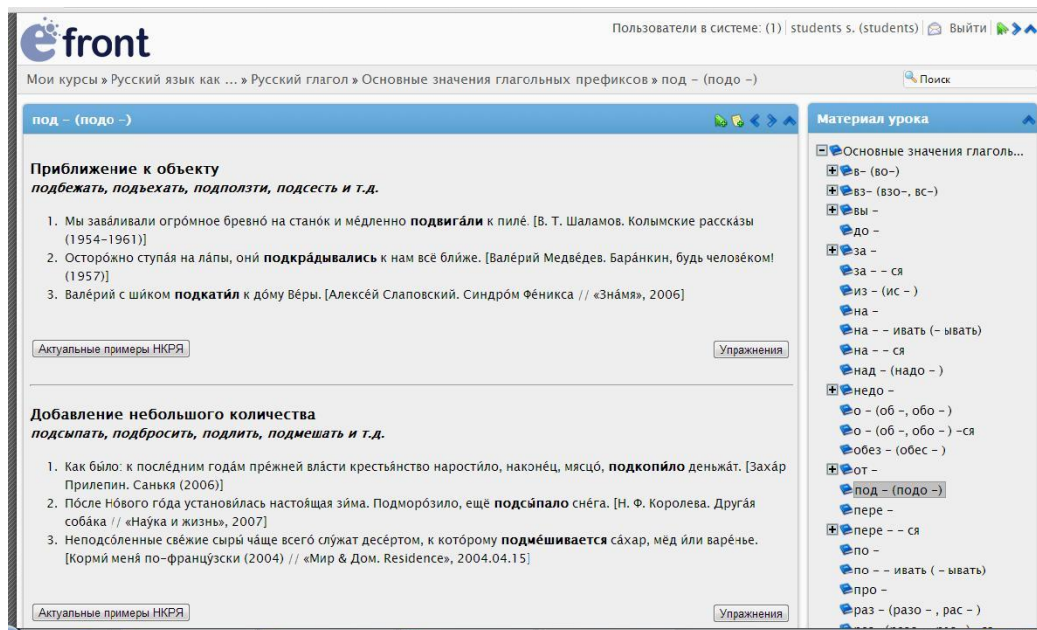


Рис. 1. Теоретический материал. Значение префиксов

Кроме того, разработана система интерактивных упражнений: созданы структурно-семантические упражнения на подстановку и множественный выбор для разных значений каждого префикса в отдельности,

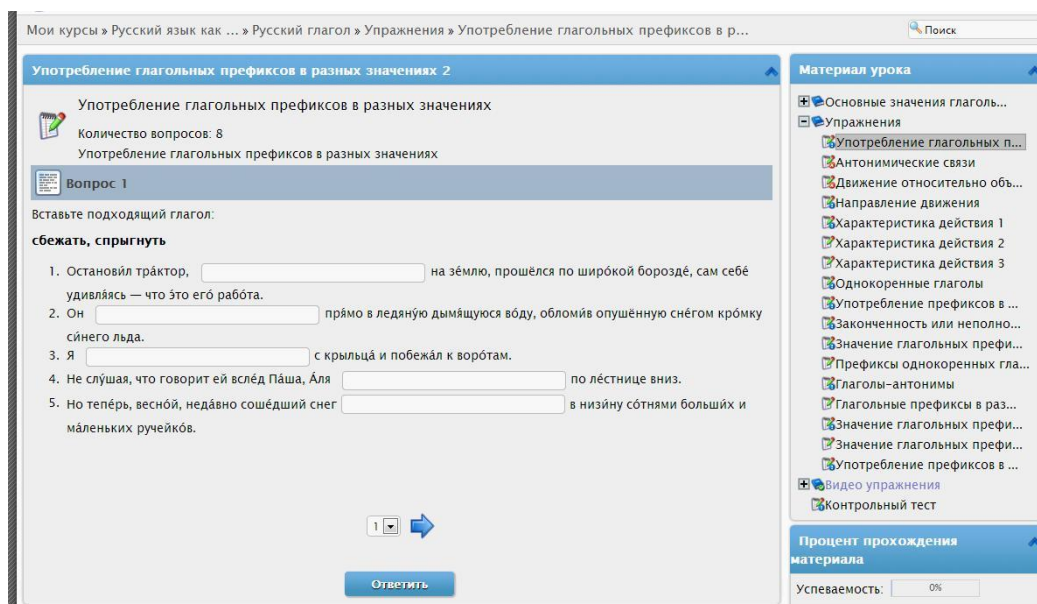



Рис. 2. Упражнения на употребление префиксов

а также придуман и реализован целый ряд упражнений, основанных на видеосоюжетах:

Вопрос 4

Какой глагол звучит в видеофрагменте?



асса_238.wmv
00:09

ископора-video, ролики пользователя

- перессориться
- пересечься
- пересекаться
- переписываться
- передрачиться

Рис. 3. Видеоупражнения

При участии иностранных студентов и преподавателей проведено экспериментальное тестирование учебного материала. Упражнения и теоретический материал размещены в системе дистанционного обучения e-front, где для удобства обучающихся собираются все материалы по русскому языку как иностранному, создаваемые научно-учебной группой «Корплинги»

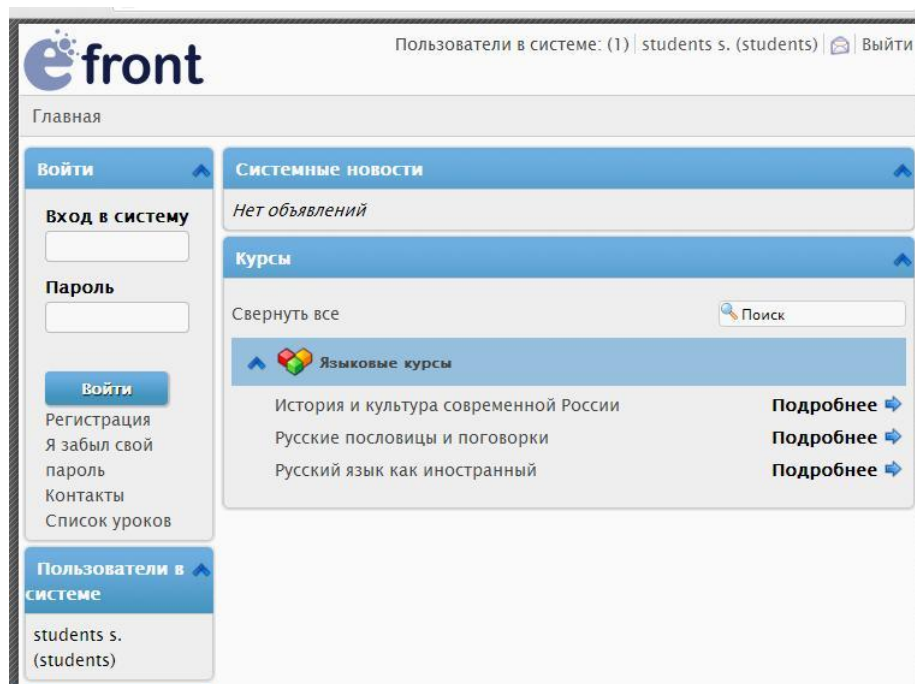


Рис. 4. e-front

Система дистанционного обучения eFront сочетает в себе функции систем управления обучением (LMS - Learning Management System) и систем управления и создания учебных материалов (LCMS - Learning Content Management System). Основу составляет eFront Core - платформа, распространяемая со свободной лицензией, реализующая основные функции LMS/LCMS. eFront имеет поддержку формата SCORM, что особенно важно для создаваемых упражнений и интеграции работы с НКРЯ. Удобство использования системы eFront заключается в возможности он-лайн работы и дистанционного выполнения заданий, это расширяет круг пользователей создаваемого учебника по русскому языку [1].

При создании учебного пособия авторы исходили из того, что современный электронный учебник не должен быть только сборником слов и правил, его цель – научить использовать язык в качестве реального (а не имитационно-подстановочного) средства коммуникации и основного инструмента научного исследования. Для существующих электронных учебников характерны недостатки, которые предполагалось преодолеть при создании электронного пособия «Русский глагол»: организация материала по принципу «от формы к содержанию»; отсутствие актуальных текстов в целях демонстрирования контекста употребления и функционирования того или иного языкового явления; незначительное количество упражнений, формирующих коммуникативную компетенцию; заметная нехватка речевых и коммуникативно-направленных упражнений. Кроме того, основное иллюстративное содержание, предлагаемое в традиционных учебниках на бумажных и электронных носителях, не меняется десятилетиями, а для преодоления постепенно увеличивающегося разрыва между сведениями, получаемыми в реальной жизни, и почерпнутыми из учебников, был необходим приток нового материала, охватывающего все функциональные стили языка. Создание электронного учебника с использованием НКРЯ было связано с «обращением к реальным современным текстам, с которыми учащийся сталкивается в своей повседневной жизни. Именно такие тексты ему придется создавать и анализировать в быту и в профессиональной сфере. Возможность быстрого получения статистических данных с помощью Корпуса позволяет составить представление о частотности того или иного явления и,

на основании этого, отобрать те факты языка, которые должны быть учтены образовательным процессом в первую очередь» (цитируется по материалам образовательного портала НКРЯ: <http://studiorum.ruscorpora.ru/>).

Благодаря продуманной метаразметке работа с НКРЯ позволила сократить время, затраченное на подбор иллюстративного материала, но в то же время потребовала поиска оптимального пути к автоматизированной выдаче результатов, поэтому работа участников научно-учебной группы вышла за рамки примитивного механического сбора примеров и позволила вести перспективное научное исследование. Методы и концепции, лежащие в основе созданного пилотного проекта учебного пособия, отличаются принципиальной новизной по сравнению с традиционными учебниками по русскому языку как иностранному: во-первых, в ходе работы было принято решение сосредоточиться только на наиболее частотных явлениях языка, которые составляют основу языковой компетенции. Во-вторых, в пособии использовались методы интерактивного обучения и регулирования, что предполагало индивидуальный подбор заданий, профилирование материала в соответствии с уровнем и потребностями отдельных студентов: учебное пособие построено таким образом, что переход от теоретического материала к упражнениям возможен на любом этапе обучения. Нет необходимости просматривать весь теоретический материал целиком, чтобы затем перейти к тренировке в упражнениях. Каждый теоретический фрагмент имеет «кнопку» перехода к упражнениям и всегда есть возможность вернуться к теории. Это способствует работе в индивидуальном режиме. Теоретическая часть расширена «Актуальными примерами НКРЯ», которые семантически соответствуют тематике, но извлекаются из Корпуса в случайном порядке.

Наделение признаком
обогащать, округлить, ослепить, обновить и т.д.

1. Поса́дку в каби́ны за́мётно **облегча́ют** внутренне́е ру́чки на пе́редних сто́йках. [Анато́лий Карпенков, Ю́рий Нечетов. Балти́йские голова́стики (2003) // «За руле́м», 2003.05.15]
2. Автор счита́ет, что да́нное напра́вление а́нализа сле́дует продо́лжить, что **обогати́т** на́ши зна́ния о потре́бительских за́пасах. [Потре́бительские за́пасы – су́щность и подхо́д к а́налізу // «Вопро́сы ста́тистики», 2004]
3. Ещѐ египтя́не **освежа́ли** во́ду с по́мощью ме́ди. [Ме́дь протѝв бакте́рий // «Зна́ние – си́ла», 2003]

Актуальные примеры НКРЯ Упражнения

1. Уже в середине ночи мы **обновили костер**, и я сделал из шоколадной фольговой обертки бокал для Ирены. [Константин Воробьев. Вот пришел великан (1971)]
2. Постсредневековая Европа **обогастила лексикон** таким понятием, как Реформация. [Алексей Муравьев. Компромисс как тактика и стратегия победы // «Отечественные записки», 2003]
3. Подкрасил кое-что, **обновил систему** отопления и вентиляции, поменял старую входную дверь все своими руками. [Руслан Хестанов. Как пережить кризис // «Русский репортер», № 3 (33), 31 января – 7 февраля 2008, 2008]
4. Советские ученые **обогастили науку** новыми, важными сведениями, заполнили пробелы географических карт. [Б. Галин. Над картой Родины // Культурная жизнь, 1947]
5. Махинации с выводом войск из Восточной Европы, Прибалтики, с дележом армии между бывшими республиками сказочно **озолотили генералов**, превратив их в крупнейших феодалов-богачей, кому на «окоормление» были сданы рода войск округа и армии с бесплатной солдатской рабсиллой. [Владислав Шурыгин. Марш побежденных (2003) // «Завтра», 2003.01.01]
6. Пусть сойдет снова Христос и **обновит души**, а иначе в человеке все порядочное исчахнет и издохнет от смрада ваших материальных благ. [А. Ф. Писемский. Мещане (1877)]
7. Находка их **обогатит человечество** неизвестными философскими системами и перлами поэзии. [К. Г. Паустовский. Повесть о жизни. Начало неведомого века (1956)]
8. Доселе государство Василия было славно и счастливо: он усилил великое княжение знаменитыми приобретениями без всякого кровопролития; видел спокойствие, благоустройство, избыток граждан в областях своих; **обогатил казну** доходами; уже не делился ими с Ордою и мог считать себя независимым. [Н. М. Карамзин. История государства Российского: Том 5 (1809–1820)]

Рис. 5. Актуальные примеры

База данных для учебников по РКИ

Следующий этап работы научно-учебной группы связан с адаптацией языкового материала НКРЯ для базы данных электронного учебника «Русский язык как иностранный». Методы сбора, обработки и классификации языковых примеров, а также разработка компьютерной оболочки для представления продукта в интерактивном, дружественном к пользователю виде, освоенные в ходе работы над проектом «Русский глагол», стали закономерным этапом перехода от научно-прикладных задач к экспериментально-научным.

Идея адаптации материала НКРЯ родилась не случайно: несмотря на то, что иностранными студентами, изучающими русский язык, а также экспертами и преподавателями РКИ высоко оценен новаторский характер пилотного пособия «Русский глагол», сложность «живого» материала подчеркивалась неоднократно. Для создания алгоритма адаптации контрольные тексты «упрощаются» двумя независимыми преподавателями-экспертами, которые описывают шаги адаптации текстов (разбивка сложных предложений на простые, замена малоупотребимых синтаксических структур, замена синонимами лексических единиц, не соответствующих выбранному уровню обучения) и систематизируют их, создавая алгоритмическое предписание. Эффективность работы интеллектуальных обучающих систем, основанных на лингвистических правилах, подчеркивается в работах многих исследователей [2, 3]

Разработчиками принято решение адаптировать языковой материал при помощи программных средств (как на уровне лексики, так и на уровне синтаксических структур) и экспериментально выявить наиболее благоприятный способ адаптации: синтаксический, морфологический или наложение этих способов.

В первом случае интересно исследовать синтаксические структуры предложений, на которых построено обучение русскому языку в соответствующих учебных пособиях (условно можно обозначить такие предложения как «простые»). В то же время, многие предложения, хранящиеся в Корпусе и представляющие реальный узус русского языка, существенно превосходят примеры из учебников по синтаксической сложности (например, содержат несколько вложений подчиненных клауз, распространенные причастные и деепричастные обороты, вводные подчиняющие обороты и т.п.). Такие предложения можно назвать «сложными». Во втором случае представляется логичным сопоставить лексическую сложность примеров в традиционных учебниках и корпусных примерах: если в первых количество незнакомых лексем не превышает 7-15% для того или иного уровня владения языком, то последние могут включать в себя малоупотребимые слова и специальную лексику. Как правило, если объем незнакомых пользователю слов превышает критический уровень, это существенно затрудняет выполнение учебных заданий (к примеру, по грамматике).

В целом задачу второго этапа проекта можно было сформулировать как автоматическую идентификацию лексической и синтаксической сложности данных НКРЯ, а также автоматическое преобразование «сложных» предложений в «простые» в целях создания базы адаптированных примеров для учебника по РКИ. Выявление правил преобразования и адаптации синтаксической структуры, чтобы при этом не терялся смысл, может дать импульс теоретическим исследованиям в когнитивной области. Верификация степени адаптации аутентичного речевого материала без искажения смысла в полуавтоматическом режиме послужит привлечению интереса к этой области работы с языком.

Программа идентификации лексической сложности успешно реализована; она имеет базовый словарь лексического минимума 1 сертификационного уровня и пополняющийся словарь однокоренных структур.

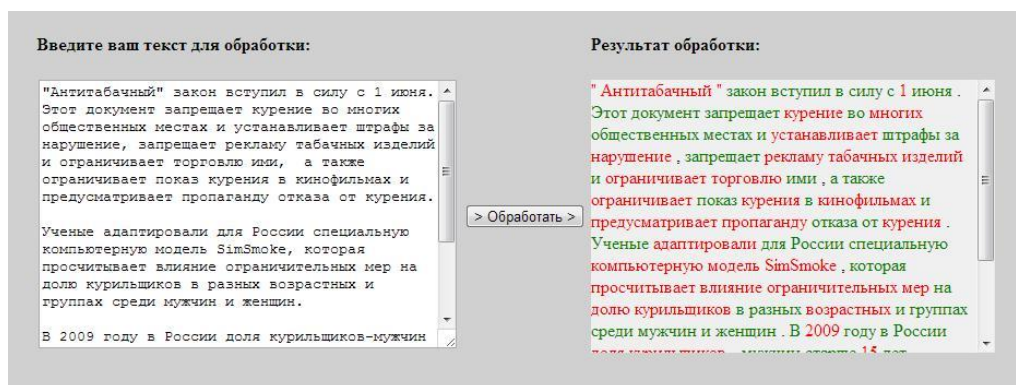


Рис. 6. Обработка лексики

При работе на втором этапе (пополнение материалов электронного учебного пособия), коллектив научно-учебной группы попытался определить способ, по которому будет строиться автоматический алгоритм: нужно было в том числе выбрать, будет ли это путь отсечения сложных для восприятия, не входящих в грамматический минимум конструкций либо это будет алгоритм, основанный на грамматическом подобии (имеем образцовые, достаточно простые грамматические структуры, по их образу и подобию автоматически отбираем схожие). Было принято решение интегрировать два возможных пути, чтобы добиться наилучших результатов.

Лингвистические задачи по упрощению предложений

В задачу лингвистов входила необходимость создания коллекции правил, которые в дальнейшем могли быть запрограммированы и использованы в качестве основных:

а) для извлечения из текстового материала различной сложности наиболее простых предложений путём отсечения сложных для восприятия, не входящих в грамматический минимум конструкций;

б) для упрощения сложных для восприятия, не входящих в грамматический минимум конструкций.

Во-первых, необходимо было изучить актуальные грамматические справочники и грамматические минимумы для обучения русскому языку как иностранному [4, 5, 6]. Во-вторых, потребовался анализ грамматической системы русского языка в соответствии с нормативными справочниками и пособиями, а также курсами лекций по морфологии и синтаксису русского языка. Путём сравнения материала, который в соответствии со стандартами для изучения РКИ на базовом и первом уровнях должен присутствовать в грамматическом минимуме, с обобщённой грамматической системой русского языка вручную были выделены структуры (синтаксические, семантические усложнители, а также сложные предложения), которые слишком трудны для восприятия на базовом и первом уровнях владения русским языком.

Sentence	ID	Class	Comment
Большую роль в развитии армянского сегмента в интернет играют международные организации и ...	[202]		simp
13 июня 2003 года открылся проект НАТО "Виртуальный шелковый путь".	[203]		simp
Только в Ереванском научно-исследовательском институте механики работало десять тысяч сот...	[204]		simp
Он производил аппаратные средства ЭВМ и вел разработку программного обеспечения.	[205]		? аббревиатура ЭВМ
Экспорт продукции данной отрасли в 2000 году составил 20 миллионов долларов США.	[206]		simp
Число ИТ-компаний в стране растет.	[207]		? сокращение "ИТ"
Так, благодаря частным и иностранным инвестициям начали создаваться специальные структуры...	[208]		! более 8 членов
Сетевые и компьютерные технологии востребованы в стране.	[209]		simp
Так, с конца 2002 года Армения начала практиковать услугу получения виз через интернет.	[210]		simp
Сегодня основными рынками экспорта для армянских ИТ- предприятий являются США, Россия, Ге...	[211]		! более 8 членов
Перспективы развития сферы высоких технологий	[212]		simp
Армения стала седьмой постсоветской республикой-членом ВТО.	[213]		! сложное слово "ре...
Доля ВВП на душу населения в 2001 году составила 3,350 долларов США.	[214]		simp
Объем ИТ-рынка растет.	[215]		? сокращение "ИТ"
Недавний визит президента РА в Москву стал важной вехой в жизни Армении.	[216]		? "РА, вехой"
Тому уже есть достаточно много примеров.	[217]		simp
В 2003 году компания "Гранд Холдинг" начала оказывать финансовое содействие в строительстве...	[218]		! более 8 членов
В Армении сегодня примеров подобного взаимодействия знаний и производства крайне мало.	[219]		simp
В большей степени имеющиеся производства и учреждения не компьютеризированы.	[220]		simp
Средств на приобретение современной техники и оборудования явно не хватает.	[221]		simp
АРТИСТ МИМАНСА	[222]		!
Это случилось вечером, в третьем акте.	[223]		simp
По зрительному залу пошел смех	[224]		simp

Рис. 7. Простые и сложные предложения

Осложнители в соответствии с синтаксической системой русского языка были классифицированы следующим образом:

- коммуникативные: модальность, эмоциональная окрашенность (междометия, вводные компоненты, обращения);
- структурные: любые обороты, присоединительные конструкции, вставные конструкции;
- структурно-семантические осложнители.

Важно, что осложняющий компонент может быть выражен со структурной точки зрения любой языковой единицей: отдельной словоформой (чаще всего акцентируется с помощью частиц, союзов и др.); сочинительным рядом; словосочетанием; оборотом (грамматической конструкции из главного слова и зависимых слов); предложением.

Сложные предложения были классифицированы с точки зрения грамматических связей и формальных показателей: наличие нескольких глаголов одной грамматической формы (грамматический показатель), характерные для разных видов связей союзы (формальный показатель) и пр.

После создания классификации была проделана работа по заполнению слотов классификации, то есть непосредственно работа по написанию правил, позволяющих отсекалть сложные, не входящие в грамматический минимум конструкции.

На данном этапе возникла проблема, связанная с отсутствием удовлетворяющего потребностям работы синтаксического анализатора, который в автоматическом режиме вычленил бы осложняющие элементы. Эта проблема была решена следующим образом: вся синтаксическая классификация осложнителей и иных элементов, делающих ту или иную структуру не подходящей для грамматического минимума, была переведена на язык морфологии (по сути, были составлены некие структурные схемы, которые были понятны компьютеру). Это было сделано, чтобы имеющаяся в нашем распоряжении программа, написанная на языке Python, объектно-ориентированном, интерпретируемом, переносимом языке сверхвысокого уровня, с использованием библиотеки PyMorphu (библиотека для морфологического анализа) могла обрабатывать материал. Программирование на Python позволяет получать быстро и качественно необходимые программные модули. Интерпретатор Python может быть перенесён на любую платформу, будь то Unix, Windows, Linux, RiscOS, MAC, Sun [7, 8].

Для корректной работы программы необходимо было также учитывать и некоторые формальные элементы, например, наличие пробелов, расстановку знаков препинания.

В конечном итоге получили лингвистические правила следующего вида, например, для выделения простых и сложных предложений:

Предложениями считать группы между .;:?!... и запятыми, если между двумя запятыми стоит три и более слов. Слово после .?!... должно начинаться с большой буквы после пробела. Группы из 2-3 слов сортируются отдельно и вручную определяется полнота структуры.

На этапе вычленения усложнителей появилась проблема выделения, классификации и преобразования неполных, эллиптических конструкций (например, взятых из диалогической речи). Одним из способов решения этих проблем может стать совмещение работы синтаксического анализатора и ручной семантической разметки, которые бы координировали элементы синтаксической и смысловой структур предложений.

Основной задачей на данный момент является описание коллекции правил, входящих в синтаксический минимум по РКИ первого сертификационного уровня, а также непрерывное пополнение коллекции так называемых «запрещающих» правил (описывающих структуры, которых не должно быть в базе адаптированных предложений). Необходима их интеграция с правилами, описывающими простые структуры, входящие в синтаксический минимум, в разрезе их наложения друг на друга и улучшения работы программы.

Решение задач технического характера

Одна из важных задач, поставленных перед программистами – участниками научно-учебной группы – автоматическое определение синтаксической сложности предложений в тексте.

Для решения данной задачи могут использоваться два подхода: машинное обучение, основанное на выборке из уже размеченных по категориям сложности предложений, и использование экспертных правил, составленных лингвистами на основе морфологической структуры предложений. В дальнейшем будет использована комбинация обоих подходов.

Обучение классификатора строится на основе встроенных средств библиотеки NLTK (Natural Language Toolkit) для языка Python 2.7. На данный момент реализация этого метода в NLTK поддерживает только бинарную классификацию, но этого вполне достаточно для решения поставленной задачи, так как необходимо разделить входных данных только на два класса.

В качестве обучающей выборки используется морфологически и синтаксически размеченная часть НКРЯ. Каждое предложение имеет тэг структурной сложности (simp – для простых предложений, comp - для сложных), который был присвоен ему в результате ручной разметки. Общий объём размеченного корпуса – около 4000 предложений, 80% которых будут использоваться для обучения классификатора и 20% для проверки качества обучения.

За единицу обучающей выборки принимается размеченное предложение, тэги которого преобразованы в вектор именованных признаков. Признак содержит в себе две составляющие – имя признака и значение признака. Каждая запись с набором признаков предложения сопровождается пометкой о принадлежности к одной из двух категорий сложности. Так как размеченный корпус хранится в виде XML-документа, извлечение нужных для обучения данных производится с помощью встроенных средств Python для обработки XML.

Одной из последующих задач является предобработка извлечённого набора векторов признаков – выбор способа представления обобщённых признаков

предложения и нормализация полученных данных. Это непосредственно влияет на будущее качество обучения модели. На данный момент планируется реализовать два способа представления признаков предложения: на основе морфологии и на основе синтаксиса. Каждое слово в предложении сопровождается строковой записью морфологических характеристик (вида «S ЕД СРЕД ИМ НЕОД»: существительное, единственного числа, среднего рода, в именительном падеже, неодушевленное) и синтаксической характеристики (типом связи, например, «предик»: предикат). В качестве условных обозначений в программе используются аналогичные условным обозначениям, применяемым в Открытом корпусе русского языка (<http://opencorpora.org/>).

Использовать в качестве признаков морфологические характеристики может быть нецелесообразно из-за недостаточно большого объема обучающей выборки и слишком большой вариативности морфологических признаков. Более продуктивным представляется использование синтаксических связей. В таком случае нормализация заключается в том, что каждая обучающая запись представлена в виде вектора длиной в количество всех возможных типов синтаксических связей в качестве имён признаков и количеством этих связей в предложении в качестве значения каждого признака.

Дальнейшая работа может показать, что при использовании нормализованных векторов синтаксических признаков метод опорных векторов окажется не самым продуктивным, поэтому по мере накопления материала будет проводиться исследование качества обучения с использованием различных способов классификации. Следует заметить, что выбор методологии составления выборки и выбор метода обучения реализуется уже на данном этапе, но при использовании синтаксической разметки в дальнейшем данная задача будет связана с обучением синтаксического анализатора. В качестве базовой модели используется доработанный malt-парсер, который будет обучен на той же выборке, что используется для задачи классификации.

Перспективы разработки комплексного электронного лексико-семантического пособия с интегрированной базой учебных текстов и возможностью их адаптации

Поставленная на первом этапе работы задача создания электронного пособия по русскому как иностранному на базе НКРЯ была успешно реализована. Благодаря использованию Корпуса материал пособия максимально приближен к реальному узусу и имеет новаторский характер. Пособие предназначено для применения на продвинутых этапах преподавания русского языка как иностранного (далее – РКИ), в том числе при дистанционном обучении. Учебный текстовый материал подготовлен на основе актуальных и обширных корпусных данных, при минимальной адаптации последних, пособие снабжено тренировочными упражнениями и контрольным тестом.

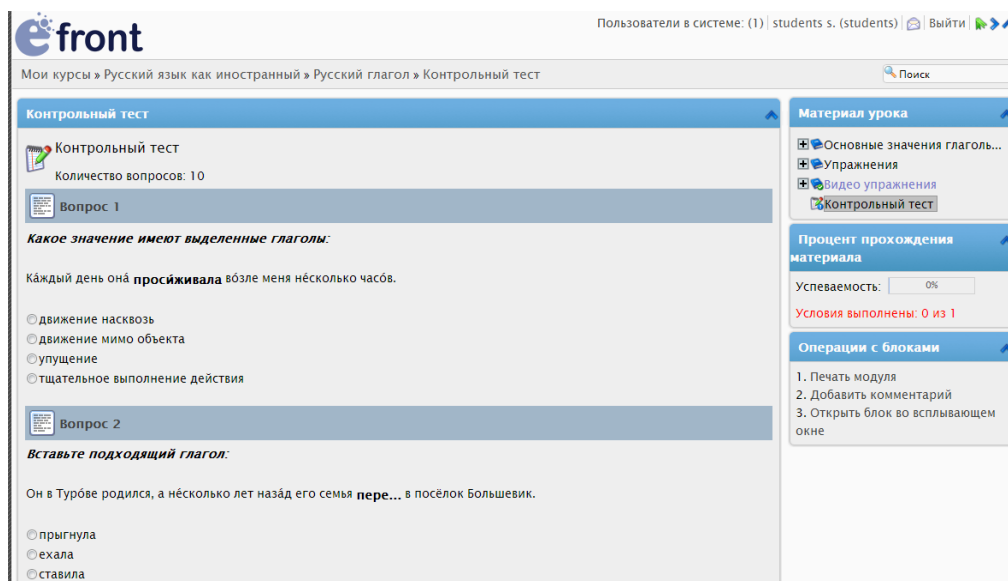


Рис. 8. Контрольный тест

Однако следует подчеркнуть, что пособие имеет очень ограниченный спектр: рассмотрена только проблематика продуктивного префиксального словообразования, включая семантику и поведение приставочных глаголов. Целый ряд важных аспектов употребления глагола, которые нуждаются в специальной отработке, остался за пределами исследования. При этом необходимо учесть сложность и многомерность такого феномена, как русский глагол, принять во внимание огромный объем и разнообразие материалов НКРЯ, а также тот факт, что проект, имеющий пилотный характер, изначально не претендовал на полноту охвата теоретических свойств и текстового массива. В перспективе возможно создание полноценного учебника по семантике и функционированию русского глагола, а в дальнейшем – и других частей речи. Но работа над данным этапом проекта считается завершенной, поскольку важнейшей задачей научно-учебной группы является поиск ноу-хау и алгоритмов создания учебного пособия по русскому (и не только) языку как иностранному, а не разработка всеобъемлющего электронного учебника.

Созданный лексический анализатор позволяет оценить сложность текстов, используемых для обучения и скорректировать употребление лексики. Немаловажным достоинством является работа с анализатором в он-лайн режиме и свободный доступ к ресурсу. Синтаксические и морфологические анализаторы применяются во многих мировых языках, а значит, поиск алгоритма и опыт написания программного обеспечения для упрощения и адаптации русскоязычного текстового материала различной сложности может быть в дальнейшем использован в работе с другими мировыми языками.

В дальнейшем разработанные технологии могут применяться и при создании блоков, представляющих другие морфологические классы и другие уровни языка (синтаксический, лексический). Созданный электронный учебник может использоваться и для занятий в русскоязычной аудитории, например, в процессе занятий по стилистике русского языка в рамках бакалавриата; созданный банк текстов также может быть полезным ученикам старших классов средней школы, помогать им в овладении литературной нормой, узусом, развивать языковую интуицию. Однако научная составляющая подобной работы будет иметь несравнимо меньшую ценность, и потому расширение учебного пособия возможно только при наличии конкретной заинтересованности со стороны преподавателей РКИ и стилистики русского языка.

Заключение

В процессе работы над электронным ресурсом по русскому языку как иностранному участники научно-учебной группы НИУ ВШЭ «Корплинги (Нижний Новгород – Москва)» нашли оригинальные решения для конкретных, локальных задач, как в области лингвистики, так и в сфере программирования: например, отбор большого массива примеров по заданной тематике, расстановка ударений, вывод неадаптированных примеров употребления глаголов в теоретической части учебника. Учебное пособие обрело дружелюбный и несложный для малоподготовленного пользователя интерфейс. Впервые представлена система видеопражнений, построенных на основе мультимедийного подкорпуса НКРЯ, а также лексический анализатор, позволяющий определить сложность учебных текстов.

Идея создания лексического анализатора, банка адаптированных примеров НКРЯ и программного продукта, позволяющего упрощать лексически и синтаксически сложные конструкции, демонстрирует оригинальный и в то же время универсальный подход в решении накопления современного текстового материала для обучения иностранному языку.

Деятельность научно-учебной группы «Корплинги» носит междисциплинарный и межкультурный характер: вокруг Национального корпуса русского языка намечено определенное направление на стыке теоретической лингвистики, корпусной лингвистики, лингводидактики и информатики, которое может развиваться и в рамках международного сотрудничества.

В данной научной работе использованы результаты, полученные в ходе выполнения проекта «Адаптация языкового материала НКРЯ для электронного учебника "Русский язык как иностранный"», выполненного в рамках Программы «Научный фонд НИУ ВШЭ» в 2013 году, грант № 13-05-0031.

Литература

1. URL: <http://www.efrontlearning.net/>
2. Галеев И.Х. Модель управления процессом обучения в ИОС // Международный электронный журнал "Образовательные технологии и общество (Educational Technology & Society)" - 2010. - V.13. - №3. - С.285-292. - ISSN 1436-4522. URL: <http://ifets.ieee.org/russian/periodical/journal.html>
3. Кривицкий Б.Х. Учебные электронные средства в вузе. Учебное пособие для преподавателей, повышающих квалификацию в МГУ. М.: МГУ, 2013. – 208 с. URL: <http://www.psy.msu.ru/people/krivitsky/krivitsky2013.pdf>
4. Требования по русскому языку как иностранному. Первый уровень. Общее владение – Москва – Санкт-Петербург: «Златоуст», 2007. – 89 с.
5. Государственный стандарт по русскому языку как иностранному. Базовый уровень - Москва – Санкт-Петербург: «Златоуст», 2001. – 112 с.
6. Акишина А.А., Каган О.Е. Учимся учить. Для преподавателя русского языка как иностранного – 2-е изд., испр. и доп. – М.: Рус.яз.Курсы, 2002. – 256 с.
7. URL: http://www.opennet.ru/docs/RUS/python/python_b.html
8. URL: <http://python.su/>