

Е.Г. – Интервьюируемая в мини-опросе (см. раздел 3), ж., 33 года, менеджер по персоналу большой американской международной компании в Москве

Е.И. – Интервьюируемая в мини-опросе (см. раздел 3), ж., 32 г., увольнение состоялось в 1999г., когда ей было 19 лет, в настоящее время менеджер в международном рекламном агентстве в Москве.

Е.П. – Интервьюируемая в мини-опросе (см. раздел 3), ж., 30 лет, руководитель службы персонала фабрики в Москве.

ИС 4 – Интервью автора <Р.Р.> с директором по персоналу крупной международной компании (ж., <Т.Я.>). М., март 2008 г.; продолжительность интервью: 70 мин.

ИС 9 – Интервью автора <Р.Р.> с генеральным директором мебельной фабрики(м., <В.К.>). Участвовали в разговоре также: педагог (м.) и научный сотрудник (ж.). Город С, май 2008 г. Продолжительность интервью: 100 мин.

СБ 9 – Собеседование с 9 претендентами на место директора филиала рекламного агентства; собеседование состоялось в офисе кадрового агентства в городе С в мае 2008 г. Участвуют: руководитель рекламного агентства (м., <Э.>), руководитель кадрового агентства (ж., <Е.>), секретарь кадрового агентства (ж., <О.>), кандидаты (все м.; <К1>, <К2>, ... <К9>), австрийская профессор (ж., <Р.Р.>). Продолжительность записи: 130 мин.

СБ 10 – Собеседование с тремя претендентками на место менеджера по обучению персонала; собеседование состоялось в городе С в мае 2008 г. Участвуют: директор по персоналу фирмы работодателя (ж., <С.К.>); ее заместитель (м., <А.Ш.>), менеджер по обучению персонала (ж., <С.>), менеджер по корпоративной культуре (ж., <Л.>), претендентки (ж., <К1>, <К2>, <К3>), австрийская профессор (ж., <Р.Р.>). Продолжительность записи: 50 мин.

СБ 13 – Собеседование с претендентом на место инженера; собеседование состоялось в кадровом агентстве в городе С в мае 2008 г. Участвуют: менеджер по подбору персонала (ж., <М.>), кандидат (м., Р.>), австрийская профессор (ж., <Р.Р.>). Продолжительность записи: 20 мин.

Wiener Slawistischer Almanach, Sonderband 83 (2013), 89 - 111

Бранко Тошович

ГРАЛИС-КОРПУС

0. В Университете им. Карла и Франца в Граце (Австрия; Karl-Franzens-Universität Graz) в 2007 году создан комплексный и многоязычный текстуальный и устный корпус для изучения славянских языков и обучение им.¹ Корпус в настоящее время используется славистами из различных стран в рамках коллективных и индивидуальных проектов и в процессе обучения.² Его можно использовать как многоязычный (параллельный), а

The screenshot shows the homepage of the GRALIS Corpus. At the top, there's a header with the university logo, the name 'GRALIS', and a search bar. Below the header, there's a navigation menu with links like 'Gralis', 'Projekte', 'Gralis-Online-Projekte', 'Gralis-Korpus', 'Lehrveranstaltungen', and 'Gastvorträge'. A timestamp 'Last Published: 01/04/2013 09:54:21' is also visible. The main content area is titled 'Das Gralis-Korpus' and contains sections for 'Das Gralis Text-Korpus' (with links to 'Das Gralis-Korpus Team', 'Die Struktur des Gralis Text-Korpus', and 'Entwicklung des Korpus'), 'Das Gralis Text-Korpus' (with a detailed description of the corpus's purpose and development), and a search interface for the corpus itself.

также как одноязычный корпус, но в первую очередь предназначен для параллелизации языков генетически близких (славянских со славянскими) и неродственных (славянских с германскими). Типологически Гралис-

¹ <http://www-gewi.kfunigraz.ac.at/gralis/index.html>.

² Руководителем проекта является Бранко Тошович (Branko Tošović), координатором работ – Арно Вониш (Arno Wonisch), работу над развитием программной поддержки возглавляет Хуберт Штиглер (Hubert Stigler), а администратором корпуса является Ольга Ленер (Olga Lehner).

Корпус относится к аннотированным корпусам (морфосинтаксическая разметка пока сделана для сербского, хорватского, боснийского и черногорского языков). Пользователи имеют в распоряжении оболочку (интерфейс) на английском, немецком, русском и сербском/хорватском/боснийском (число языков постепенно увеличивается).

Параллельные корпусы

1. Существует мало параллельных корпусов для славянских языков (далее СЛЯз). Их нет так много и для неславянских языков. Среди русско-неславянских корпусов выделяются те, которые разработаны в рамках „Национального корпуса русского языка“ и которые охватывают восемь языков: английский, белорусский, испанский, итальянский, немецкий, польский, украинский, французский (*Ruscorpora-www*),³ среди которых один – Русско-немецкий и немецко-русский корпус (автор Сергей Шаров) размещен в Лидсе (*Leeds Corpus*). В рамках проекта „Opus“ (an open source parallel corpus) подготовлен набор свободно распространяемых параллельных текстов (техническая документация, корпус субтитров) на пяти языках (*Opus-www*). Он основан на CWB CQP. В рамках проекта „Корпус Австрийской Академии наук“ готовится „Русско-немецкий корпус параллельных текстов“, охватывающий один текст – роман Ф. М. Достоевского „Идиот“ (1868–1869) и его переводы на немецкий язык.⁴ В процессе создания находится немецко-русский корпус „Толкование снов Фрейда“ – „Traumdeutung“ (*Traumdeutung-www*).⁵

В 2001 году в Самарском университете был организован проект на тему „Корpusная лингвистика и новые информационные технологии“ в рамках межкафедрального сотрудничества с Институтом немецкого языкоznания Университета г. Вюрцбурга (руководитель Н. Р. Вольф). Целью проекта является последовательное составление и компьютерная обработка параллельного немецко-русского корпуса переводных текстов (*DER-Korpus*) для исследовательских и учебных целей. Сюда относится и Параллельный корпус переводов „Слова о полку Игореве“ (*SPI-www*). Здесь можно упомянуть и „Lilabar“ – англо-русский корпус параллельных предложений, пословиц и фраз (*Lilabar-www*).⁶ На Отделении переводов Тампер-

³ Полная информация об интернет-источнике находится в списке литературы в конце работы.

⁴ Этот корпус ориентирован на изучение лексической семантики в сопоставительном аспекте.

⁵ Целью данного проекта является создание и морфологическое аннотирование немецко-русского корпуса для различных пользователей.

⁶ Он содержит базу параллельных предложений („память переводов“, англ. *translation memory*) с возможностью поиска по ней. В базе представлены переводы английских и

ского университета (Финляндия) создан русско-финский параллельный корпус художественных текстов „ПарРус“.⁷ К отдельной группе относится параллельный корпус, состоящий, с одной стороны, из текстов на английском и немецком языках, а с другой – из текстов на славянских языках (Bg, Be, Bs, Cs, Hr, Pl, Ru, Sk, Sr, Uk). Это The Regensburg Parallel Corpus (RPC) в Институте славистики Университета в Регенсбурге (RPC-www)⁸, который переименован в ParaSol (PARAllel corpus of Slavic and Other Languages) и в настоящее время является совместным проектом Института славистики в Берне и Университета Регенсбург. Корпус сопровождает морфосинтаксическая аннотация. В этом университете разработан еще один корпус – Bilingualkorpus: Das Regensburger Korpus slavisch-deutscher Bilingualer (ReBiSlav):. В его основе находятся семинары с упражнениями (Hauptseminaren mit Übung) и одна магистерская работа. Он содержит ряд интервью со славянско-немецкими двуязычными говорящими различного уровня (сами студенты делают записи и транскрибируют их). Сюда также относится „Восточный многоязычный корпус“, охватывающий болгарский, чешский, эстонский, венгерский, румынский и словенский языки.

Славяно-славянских параллельных корпусов не так много. Здесь можно упомянуть „Русско-словацкий параллельный корпус“ и „Украино-русский параллельный текстовый корпус“, состоящий из веб-публикаций, доступных для поиска в режиме онлайн и для скачивания (*Infostream-www*).

Среди неславянских параллельных корпусов выделяется Europarl Parallel Corpus (EPC-www) – открытый корпус Европарламента на 11 языках, в рамках которого развиты подкорпусы Bulgarian-English, Czech-English, Danish-English, German-English, Greek-English, Spanish-English, Estonian-English, Finnish-English, French-English, Hungarian-English, Italian-English, Lithuanian-English, Dutch-English, Polish-English, Portuguese-English, Romanian-English, Slovak-English, Slovene-English, Swedish-English. Центр теории перевода Университета Leeds развил „Leeds Corpus“ (*Leeds-www*), охватывающий различные языки (английский⁹, арабский, греческий, итальянский, испанский, китайский, немецкий, польский, португальский,

руссских предложений (фраз) в обоих направлениях. Предложения (фразы) разделены по темам.

⁷ Для проверки данных использовались относительно небольшие по объему англо-русский и русско-английский массивы художественных текстов. В его состав входят произведения русской литературы XIX–XX вв. и их переводы на финский язык.

⁸ Он охватывает только один функциональный стиль – литературно-художественный (тексты в подлиннике и в переводе).

⁹ Английский корпус состоит из новостей агентства Reuters. Использование корпуса ограничено только в целях исследования.

русский¹⁰, финский, французский, японский). Одним из параллельных корпусов является „MAASTR“ (Maastr-www), охватывающий тексты Маастрихтского договора на немецком и английском языках. В рамках немецко-французского проекта „Коллокации в контексте“ создан параллельный корпус „Немецко-французский словарь коллокаций“, включающий немецкие тексты с французским переводом и французские тексты с немецким переводом (Kollokation-www).¹¹ Это по сути словарь сочетаемостей, охватывающих типичное и постоянное окружение какого-либо лексического элемента, в первую очередь прилагательных и существительных. К этой группе примыкает также „Parallel Corpus of Portuguese and English“, сокр. COMPARA (Compara-www).¹² В 1999 году работа началась в Институте немецкого языка в Мангейме (Institut für Deutsche Sprache, IDS) работа над проектом „GeFrePac“ (German-French Reciprocal Parallel Corpus), финансируемым ELRA (European Language Resources Agency, Paris) и IDS, под руководством Вольфганга Тойберта (Wolfgang Teubert). Сюда относятся и другие корпусы (некоторые из которых являются лишь попытками), в основном ориентированные на английский язык: „Английско-германский корпус по переводу“ – „Chemnitz German-English Translation Corpus“, Английско-чешский корпус „Kacenka“, Англо-китайский параллельный корпус „HKUST“, созданный в Гонконге, Англо-норвежский параллельный корпус, Английско-французско-испанский корпус „Lancaster’s ITU“, Англо-хорватский параллельный корпус, „Intersect“ – англо-немецкий параллельный корпус, „Agenda 21“, включающий датский, английский, французский и немецкий языки, „Многоязычный корпус“, содержащий переводы Библии на английском, вьетнамском, греческом, датском, испанском, латинском, финском, французском, шведском, языках и т. д.

2. Создание Гралис-корпуса было осуществлено в рамках (а) научно-исследовательского проекта (2006–2010), финансируемого Австрийским фондом для поддержки научных исследований (FWF: Fonds zur Förderung der wissenschaftlichen Forschung: FWF-Projekt, P19158-G03), (б) сотрудничества между Институтом славистики и Центром моделирования информации Гуманитарного факультета Университета им. Карла и Франца в Граце (Zentrum für Informationsmodellierung in den Geisteswissenschaften an der Karl-Franzens-Universität Graz) и (в) при поддержке специалистов в обла-

¹⁰ Он содержит лишь тексты новостей из „Известий“ в период 2000–2001 г. (объем 14 564 884 словоупотреблений). Здесь использован и русский „Referenzkorpus“ (50 512 584 словоупотреблений).

¹¹ Он состоит из CELEX-документов (право Европейского Содружества – соглашения, внешние отношения, законы) и документов Европейского парламента (EUROPARL).

¹² При его создании особенно подчеркивалась проблема авторских прав и указывалось на то, что тексты, авторами которых являются умершие 70 или больше лет назад, не нуждаются в получении разрешения.

сти славянской корпусной лингвистики.¹³ В качестве платформы для корпуса используется Gralis – славистический лингвистический портал Университета им. Карла и Франца в Граце)¹⁴, который в 2011 году отметил десятилетний юбилей.

В данном проекте преследуются следующие цели: создать параллельный корпус для всех Сляз, пополнять его сбалансированным материалом, отражающим языковое разнообразие, проводить лингвистический, социолингвистический и психолингвистический¹⁵ анализ данных корпуса в рамках совместных и индивидуальных проектов, а также отдельных исследований.

3. Гралис-Корпус развивается в первую очередь для того, чтобы ответить на вопрос, в чем отличаются славянские литературные языки (Сляз) между собой и насколько они далеки друг от друга. В нем представлены все Сляз (14): белорусский (Be), болгарский (Bg), боснийский/бошняцкий (Bs), лужицко-сербский (Ls),¹⁶ македонский (Mk), польский (Pl), русский (Ru), сербский (Sr), словенский (Sl), словацкий (Sk), украинский (Uk), хорватский (Hr), черногорский (Mo), чешский (Cs), а также немецкий (De). По мере необходимости будут включаться и так наз. славянские микроязыки: градишчанско-хорватский (Hg), кашубский (Ks) и русинский (Rs).¹⁷

Электронные корпусы и близость/отдаленность языков

4. О степени близости/отдаленности Сляз существуют различные мнения. Они особенно поляризованы по отношению к очень близким Сляз, какими в частности являются Bs, Hr, Mo и SR или же Bg и Mk. Противоположные взгляды особенно проявляются в социолингвистическом определении характера Сляз – являются ли они (а) отдельными языками, (б) вариантами одного (полицентрического) языка, (в) диалектами одного национального языка, (г) одним языком с различными (политическими) названиями и т. п. Отсутствие исследований, которые бы на базе электронных корпусов объективно представили соотношения Сляз и предоставили релевантный материал для разъяснений всех существенных вопросов, создает благопри-

¹³ В этих целях был проведен в летнем семестре 2006 года семинар на тему „Славянская корпусная лингвистика“.

¹⁴ <http://www-gewi.uni-graz.at/gralis>.

¹⁵ Имеется в виду восприятие языков носителями и неносителями языков, особенно близкородственных.

¹⁶ Система Гралис-корпуса позволяет вносить тексты для обоих лужицкосербских языков.

¹⁷ Латинские сокращения, указывающие на название языков, используются нами во всех работах, написанных на различных языках.

ятные условия для субъективных, тенденциозных толкований взаимодействия Сляз и политизации их межъязыковых корреляций, во что втягиваются не только языковеды, но очень часто и нефилологи (в первую очередь политики). Данная проблема сильно обострилась после распада трех славянских федеративных государств (СССР, Чехословакии и Югославии) в конце XX столетия и военными конфликтами в юго-восточной Европе в это время. Ситуация особенно осложнилась из-за различных подходов к оправданности и неоправданности официального провозглашения и кодификации некоторых Сляз в XX столетии (в середине этого века МК, а в 90-е годы Bs, Hr, Mk, Mo и Sr). На усложнение межъславянских языковых отношений повлияло и изменение статуса отдельных Сляз. Так (а) Cs и Sk (б) Be и Uk, (в) SI превратились из языков республик в языки новых государств. С другой стороны, усилились тенденции к вытеснению Ru как языка общеноционального общения в Белоруссии и Украине, в которой украинское и русскоязычное население втянуто в конфликтную ситуацию. Статус языка взаимного общения на уровне одного государства потерял и сербохорватский язык, в рамках которого после распада бывшей Югославии кодифицированы три (Bs, Hr, SR), а четвертый находится в этом процессе (Mo, провозглашенный совсем недавно – в 2007 году). Некоторые из представителей так наз. микрославянских языков требуют в новых государствах другого статуса (карпатско-русинский в Украине). Во всем этом появилась еще одна новизна – с распадом общего государства некоторые Сляз стали языками национальных меньшинств (например, Hr в Сербии или SR во Хорватии). Определенная, а порой и значительная часть русскоговорящих оказалась за пределами Российской Федерации и при расширении Европейского содружества вошла в его состав (Латвия, Литва, Эстония). Почти все спорные вопросы так или иначе связаны с процессами объединения в рамках Европейского Содружества, которому приходится считаться со сложными межславянскими языковыми отношениями и искать решения. Непростая и комплексная языковая ситуация на территории Сляз требует объективных исследований интралингвистических, социолингвистических и психолингвистических аспектов их соотношения на репрезентативном корпусном материале. Тот факт, что De находится в прямом или посредственном контакте со всеми Сляз и что немецкоязычные страны имеют разветвленное сотрудничество со славянским миром, создает необходимость включить в корпус и славяно-немецкий языковой материал. Все вышеуказанные причины побудили нас создать корпус для системолингвистического, социолингвистического и психолингвистического изучения (в первую очередь) отношений между (а) всеми Сляз и (б) Сляз и De.

На материале Гралис-Корпуса можно будет рассматривать и вопрос о том, как носителями Сляз воспринимается и оценивается расстояние между Сляз, причем в центре внимания находится взаимное понимание, а также кодовое переключение (code-switching).¹⁸ Данный корпус является важным для славистики, так как можно на однообразном и унифицированном материале проводить системное и комплексное исследование отношений между Сляз, а также их параллели с De. Результаты корпусных работ укажут на индивидуальную межъязыковую славянскую и неславянскую (славяно-германскую) близость/отдаленность. Гралис-Корпус может послужить хорошей материальной основой для дальнейшего исследования Сляз и De в их взаимных связях. Для германистики он даст полезный материал для определения расстояния между De и Сляз, языками, относящимися к различным генетическим семьям. В настоящее время в составе ЕС находится пять Сляз: Bg, CS, Pl, Sk и SI, в ближайшее время такой статус получит Hr, на что также претендуют BS, Mo, Mk и Sr, т. е. десять Сляз могут в ближайшие годы оказаться в ЕС. Одной из проблем является то, что значительная часть языков стран, рассчитывающих на членство в ЕС, относится к группе близкородственных, носители которых свободно, без переводчика общаются друг с другом (типичный пример – бошняки, сербы, хорваты, черногорцы), поэтому в рамках ЕС ведутся дискуссии о том, что делать с такими языками и надо ли для них создавать службу перевода. Создаваемый корпус может дать объективную картину о различиях между ними и о том, насколько оправдано, целесообразно это делать. Гралис-Корпус предоставит драгоценный материал тем, которые занимаются переводом, так как в его основе лежат „битексты“¹⁹ и точки

¹⁸ Некоторые исследователи считают его основным критерием для различия языков (если говорящие друг друга не понимают, они говорят на различных языках, если же понимают друг друга, речь идет об одном языке или диалекте). Хадсон приводит причины, почему к критерию взаимного понимания надо осторожно относиться, в частности: существуют варианты, которые можно считать различными языками, но их носители понимают друг друга (скандинавские языки, кроме финского и лапонского), но и у вариантов, принадлежащих к одному языку, может отсутствовать взаимное понимание (диалекты китайского языка) (Hudson 1990, 35–37). Анализируя тесты понимания как посредственного метода измерения лингвистического расстояния, Амон пришел к выводу, что отмеченную в этих тестах степень взаимного понимания нельзя считать однозначным индикатором лингвистического расстояния (Amon 1986, 13). Более подробно об измерении лингвистического расстояния и лингвистического несходства, а также о тестах по измерению взаимного (не)понимания Амон писал в другой работе (Amon 1989, 31–46).

¹⁹ Битекст представляет собой параллельный текст на одном языке вместе с его переводом на другой язык. Это название ввел Брайан Харрис (Brian Harris) в 1988 году, а концепция, базирующаяся на нем, была развита группой ученых при Университете Монреаля (Université de Montréal) под названием RALI (*Recherche appliquée en*

соприкосновения с концепцией памяти переводов.²⁰ Для студентов он особенно полезен, из-за наличия различных баз данных, вытекающих из него или связанных с ним. Кроме того, студентам предоставляется возможность собирать и исследовать материал для курсовых, дипломных, магистерских, кандидатских и докторских работ:

Когда Гралис-Корпус будет полностью готов, при его помощи можно будет искать ответ на вопрос, насколько близки/далеки языки, относящиеся к генетически одной языковой группе (славянской), и языки, принадлежащие к различным языковым семьям (германской и славянской), а также ответить на вопросы: (1) какие лингвистические, социолингвистические и психолингвистические процессы влияют на близость/отдаленность СЛЯз, (2) насколько лингвистическое расстояние между СЛЯз существует на языковую ситуацию и психолингвистические процессы, в первую очередь на порог узнавания и понимания, (3) насколько СЛЯз приближаются или удаляются друг от друга в четырех различных группах: в группе очень близких языков (например, Bg и Mk), в группе территориально близких языков – восточнославянских, западнославянских и южнославянских (например, Be, Ru и Uk) и в группе территориально отдаленных языков (например, HR и Cs).

linguistique informatique или *Applied Research in Computational Linguistics* – «Прикладные исследования в вычислительной лингвистике». „В сфере исследований в области перевода «битекст» – это совмещенный документ, состоящий из версий соответствующего текста на исходном и целевом языках. Битексты создаются с помощью специальных компьютерных программ, которые называются «инструментами для выравнивания» (*alignment tool*) или «инструментами для битекста» (*bitext tool*), которые позволяют автоматически выравнивать оригинальную версию текста и его перевод. Подобные программы, как правило, приводят в соответствие два текста (оригинал и перевод) по каждому предложению. Собрание битекстов называется «битекстовой базой данных» или «двухязычным корпусом» и может использоваться в качестве справочника и для поиска нужных сочетаний“ (<http://ru.wikipedia.org/wiki/Битекст>).

²⁰ Различие между битекстом и памятью переводов состоит в том, что память переводов представляет собой базу данных, в которой сегменты текста (соответствующие друг другу предложения) не связаны с оригинальным контекстом (оригинальная последовательность предложений теряется): „Битекст же сохраняет изначальную последовательность предложений. Стандартным форматом для обмена базами данных памяти переводов между разными программами автоматизированного перевода является формат TMX (XML словарь, опубликованный LISA (Ассоциация отрасли локализации – Localisation Industries Association). TMX позволяет сохранять оригинальный порядок предложений. Битексты создаются в качестве справочного инструмента для консультаций специалистов-переводчиков, а не автоматизированных программ. Поэтому небольшие ошибки выравнивания или неточности, которые могут привести к сбоям в памяти переводов, для них не имеют значения“ (<http://ru.wikipedia.org/wiki/Битекст>).

Целью корпуса является предоставление разнообразного материала для изучения совпадений, сходств и различий между СЛЯз. Так как речь идет о близких языках и так как их отношение служит поводом для разнообразных спекуляций, данный корпус может послужить материальной основой для объективной оценки. Гралис-Корпус преследует цель показать, как СЛЯз функционируют и взаимодействуют на конкретном материале, в реальном контексте и на всех языковых уровнях (фонетико-фонологическом, орфоэпическом, грамматическом и стилистическом). В подготовке Гралис-Корпус делается упор на то, чтобы он был как можно более презентативным и сбалансированным. Гралис-Корпус готовится в рамках концепции языковых корреляций, изложенной в книге *Korelaciona sintaksa* (Tošović 2001) и в ряде статей, особенно в работах, в которых рассматривается теоретический вопрос о том, что такое различие как понятие и различие в языке (Tošović 2008e).²¹ В качестве теоретической основы для типологической обработки и представления текстов служит книга *Funkcionalni stilovi* (Tošović 2002). Согласно концепции, изложенной в ней, Гралис-Корпус разделен на пять функциональных стилей (литературно-художественный, публицистический, научный, официально-деловой и разговорный). При создании корпуса учитывались следующие критерии: (1) он не должен зависеть от любых внешних обстоятельств, (2) он должен быть в состоянии следить за ходом и скоростью изменений в области информационных технологий и (3) иметь возможность постоянно развиваться, совершенствоваться и дорабатываться. Так как качество любого корпуса определяется (а) глубиной и широтой аннотирования, (б) возможностями поиска и получения информации, (в) представительным (репрезентативным) характером, пропорциональностью и сбалансированностью, а также (г) свободой доступа, всему этому уделяется большое внимание. Гралис-Корпус приспособлен в значительной степени к программным мировым стандартам (TEI, XML и др.). Его развитие направлено на количественное расширение новым содержанием, качественное улучшение оболочки, более широкое и глубокое аннотирование, функциональное ускорение, улучшение поисковой системы и программное обновление.

Корпус предназначен в первую очередь для специалистов в области сопоставительного языкознания и для языковедов более широкого профиля (особенно в области общего и системного языкознания, а также социолингвистики), для интересующихся интралингвистическими, интерлингвистическими и экстралингвистическими отношениями между СЛЯз. Он является полезным для всех тех, которые сталкиваются с проблемой

²¹ Некоторые теоретические и практические аспекты отношений между близкородственными языками рассмотрены и описаны в сборнике *Die Unterschiede zwischen dem Bosniakischen/Bosniakischen, Kroatischen und Serbischen* (Tošović 2008).

взаимодействия Сляз. Его можно целесообразно использовать в обучении, особенно в вузах.

Концептуальная основа Гралис-Корпуса

5. В основе стратегии развития Гралис-Корпуса лежит наша концепция славянских корреляционных систем, состоящая из интракорреляционала, интеркорреляционала, супракорреляционала, суперкорреляционала и экстракорреляционала. Интракорреляционал представляет собой сеть отношений внутри одного языка, влияющих на межъязыковом уровне (например, изменения внутри Ru, влияющие на расстояние между Ru и Pl). Интеркорреляционал образуют очень близкие языки, языки с очень высокой степенью понимания (скажем, BS, Hr, Mo, Sr). Супракорреляционал состоит из территориально близких Сляз (восточнославянских, западнославянских, южнославянских). Суперкорреляционал охватывает языки, относящиеся к территориально различным славянским группам (напр. восточнославянским и южнославянским типа Ru ↔ Bg). Экстракорреляционал является системой отношений генетически различных языков, в данном случае славянских языков с одной стороны и немецкого с другой.

Первая корреляционная система – интракорреляционал состоит из отношений в рамках только одного языка. Предметом исследования данного типа при помощи корпусного материала являются (**a**) динамические процессы, а именно изменения в отдельных языках, влияющие на межъязыковое расстояние и ведущие к увеличению или сокращению расстояния между Сляз, к (не)восприятию нововведений носителями определенного языка, к положительному или отрицательному восприятию, к увеличению или ослабеванию понимания носителями других славянских языков, т. е. к увеличению или снижению порога узнавания и понимания, (**b**) статистические процессы, а именно как и насколько сложившиеся структурно-типологические свойства воздействуют на межъязыковые связи. Для каждого языка в интракорреляционале существует главная временная линия разграничения, чаще всего важное для языковой проблематики события: например, для Be, Ru и Uk это распад ССР (1989), для Cs и Sk – распад ЧССР (1993), для Bs, Hr, Mo и Sr – распад СФРЮ (1992), а для De – объединение Германии (1990). Если нет отчетливой точки отсчета, то можно учитывать временной ориентир для других языков данной группы (например, для Pl и Ls в группе западнославянских языков – 1993 год, для Bg в группе южнославянских языков – 1992). Для определения интракорреляционного расстояния между Сляз при помощи такого корпуса необходимо иметь, как минимум, две версии из различных периодов.

Вторая корреляционная система – интеркорреляционал охватывает очень близкие Сляз, какими являются (a) BS, Hr, Mo, Sr, (b) Bg и Mk, (b) Sk и Cs. Проведенное исследование (Tošović 2008) свидетельствует о том, что процесс их дивергенции усилился после распада бывшей Югославии, что повлияло на увеличение расстояния между BS, Hr, Mo, Sr. Этому особенно способствуют радикальные пурристические тенденции в некоторых из них и обострение общественных процессов (усиление национализма и шовинизма, развязывание военных конфликтов и т. п.) и психические факторы (напр., ненависть к другим народам и их языкам). Целью корпусных работ в данной части является не только подготовка корпусного материала для изучения отношений между Сляз внутри каждого интеркорреляционала (**A**, **B** и **C**), но и вопрос о межгрупповом расстоянии (скажем, является ли расстояние между Hr и Sr большим, чем между Bg и M или между Cs и Sk).

Супракорреляционал состоит из отношений территориально близких Сляз. В его состав входит восточно-западнославянский супракорреляционал: Be, Ru, Uk (**A**), западно-южнославянский супракорреляционал: Cs, Ls, Pl, Sk (**B**) и южнославянский супракорреляционал: Bg, Bs, Hr, Mk, Mo, Sr (**C**). Цель корпусной работы в рамках супракорреляционала – подготовка материала для определения различий между языками, входящими в состав каждого супракорреляционала в отдельности (скажем, расстояния между Ls, Pl, Sk и Cs), а также определение расстояния между языками, принадлежащими различным супракорреляционалам (здесь ставится вопрос, являются ли более близкими или более далекими языки супракорреляционала **A** по отношению к **B** или **C**, напр. Bg, Mk ↔ Cs, Sk, Cs, Pl ↔ Ru, Uk, Be, Ru ↔ Mk, Sl).

В составе суперкорреляционала находятся Сляз из различных групп – 1. восточно-западнославянской: а) Be, Ru, Uk, б) Ls, Pl, Sk, Cs, 2. восточно-южнославянской: а) Be, Ru, Uk, б) Bg, Bs, Hr, Mk, Mo, Sr, 3. западно-южнославянской: а) Cs, Ls, Pl, Sk, б) Bg, Bs, Hr, Mk, Mo, Sr. Одной из корпусных проблем в данной сети взаимодействий является вопрос, насколько интеркорреляционное изменение влияет на суперкорреляционное расстояние. Здесь будет предпринята попытка собрать корпусный материал для рассмотрения положения о том, что увеличение интеркорреляционного расстояния влияет на характер суперкорреляционного расстояния. Скажем, тяготеют ли некоторые процессы в Hr, направленные на сознательное удаление от Sr (= увеличение интеркорреляционного расстояния), к сближению Hr и Ru (= сокращение супракорреляционного расстояния).

Для определения интер-, супра- и суперкореляционного взаимодействия необходимы тексты, переведенные по возможности на все Сляз или на большинство из них.

Особую систему составляют соотношения Сляз с немецким языком в рамках экстракореляционала. Здесь в центре внимания находятся связи между этим языком и славянскими интракорреляционными языками (например, Bs, Hg, Mo, Sr), супракорреляционными языками (например, Cs, Ls, Pl, Sk) и суперкорреляционными языками (например, Sl и Uk, Pl и Be, Mk и Cs). Данная часть корпусных работ преследует цель собрать материал, чтобы ответить на вопрос, к каким Сляз стоит ближе De, а от каких он находится подальше. Здесь надо проверить гипотезу о том, что наличие прямого территориального контакта влияет на сокращение расстояния между Сляз и De. Для определения экстракорреляционного расстояния необходимо выбрать (а) славянские тексты, имеющие больше всего переводов на немецкий язык (по возможности уже включенные в корпус по определению интер-, супра- и суперкорреляционного расстояния), и (б) немецкие тексты, имеющие больше всего переводов на Сляз.

Каждая из вышеупомянутых корреляционных систем должна иметь свой подкорпус – Интра-Кор (только для одного языка), Интер-Кор (для очень близких Сляз), Супра-Кор (для языков, относящихся к одной из трех групп Сляз), Супер-Кор (для языков, относящихся к различным славянским группам), Экстра-Кор (для Сляз и немецкого), и отличаться формальной унификацией и функциональной сбалансированностью.

Для исследования межславянских языковых отношений в их корреляционных системах в корпус надо включать переводы, опубликованные в новейшее время. При этом нужно учитывать то, какой из текстов является самым переводимым на Сляз. В рамках литературно-художественного стиля основным материалом являются переводы прозаических произведений. Работа над созданием корпуса для публицистического стиля подразумевает поиск текстов с одинаковым, немодифицированным содержанием из on-line изданий типа „Southeast European Times“, „Deutsche Welle“, „Voice of America“ и т. п. Так как они покрывают лишь некоторые Сляз, для остальных приходится делать переводы с уже имеющих языковых версий. При поиске материала для научного стиля упор делается на переводы собственно научных произведений. Для анализа отношений между Сляз в рамках официально-делового стиля важнейшими являются славянские версии основных документов международных организаций (ОНН, ЕС, ЮНЕСКО и др.). Что касается разговорного стиля, уже создается SlawSpeech-Korpus, и его материал используется для написания одной докторской диссертации. В его состав войдут три подкорпуса: SlawWort-Korpus, SlawFix-Korpus и SlawFrei-Korpus. Первые два будут использованы

для фонетического и просодического анализа. В SlawWort-Korpus будут вноситься записи отдельных слов, являющихся общими для всех Сляз. Например, составляется список из двадцати слов, относящихся к различным частям речи и используемых во всех Сляз. Эти слова читаются опрошенными. Запись потом разделяется на сегменты, которые соответствуют каждому слову, а потом они параллелизуются таким образом, чтобы при нажатии на одно слово появились все произнесенные версии на всех Сляз. В SlawFix-Korpus войдут записи связанных слов в предложении. Такие тексты должны быть небольшими (не более двадцати предложений). Третий подкорпус – SlawFrei-Korpus предназначен для определения отношений между Сляз на уровне текста и стиля. Он будет охватывать спонтанно произнесенные высказывания на одну из выбранных тем, что даст возможность измерять расстояние в речи, не подлежащей внешней языковой корректуре и вмешательству цензуры. Скажем, опрошенному дается рисунок или набор рисунков с просьбой рассказать, что там нарисовано. Так как данное исследование охватывает большое число языков, такие записи должны быть короткими. В SlawWort-Korpus войдут слова, общие для Сляз и De, а в SlawFix-Korpus (а) самые простые предложения на Сляз для опрашиваемых с родным немецким языком и (б) самые простые предложения на De для опрашиваемых с родным Сляз.

В процессе анализа корпусного материала проводятся следующие процедуры: (а) берется определенный материал (словарный или текстуальный) из двух различных периодов одного и того же языка (напр., 1970–1990, 1991–2010) и фиксируются изменения, которые влияют или могут влиять на расстояние (увеличение или сокращение) между Сляз, (б) определяется характер этих изменений: почему они происходят (для более экономного и эффективного выражения, спонтанно или целенаправленно, с политическими целями, для усиления межъязыковой дивергенции или конвергенции и т. п.), являются ли они релевантными, случайными, спонтанными, запланированными, целенаправленными и др.

При заполнении корпуса материалом предусмотрено пропорциональное и уравновешенное включение текстов трех типов: оригинальных, модифицированных (адаптированных) и переведенных.

На собранном корпусном материале можно изучать влияние языковой политики, стандартизации и кодификации на сокращение или увеличение расстояния между Сляз, а также рассматривать вопрос, как носителями Сляз воспринимаются и оценивается расстояние между отдельными языками.²²

²² Здесь в центре внимания находится критерий взаимного понимания, а также влияния расстояния на кодовое переключение (code-switching).

Структура Гралис-Корпуса

6. Гралис-Корпус состоит из нескольких подкорпусов, разделенных на макрогруппы (восточнославянские, западнославянские и южнославянские языки) и микрогруппы (корпусы отдельных языков, а также индивидуальные корпусы – корпусы писателей, напр., Иво Андрича, Зорана Живковича и др.). Существует также возможность выбора двух макрогрупп (скажем, восточнославянской и южнославянской). Отдельный блок составляют славянские языки, соотнесенные с немецким языком.

The screenshot shows the Gralis Corpus search interface. It includes a sidebar with filters for Group (all), Corpus (Gralis Slav-Korpus), Primary language (Russian), Author, and Functional style. There are also filters for Target language (Belarusian, Bulgarian, Bosnian/Croatian, Serbian/Croatian, German), Width of context (One sentence), and Width of page (All, Corpora text as HTML). A map of Europe is visible in the background.

7. В рамках Гралис-Корпуса выделяются две основные части: Корпус устной речи (Speech-Korpus) и Корпус письменных текстов (Text-Korpus).

8. Устный Гралис-Корпус (Gralis Speech-Korpus) охватывает транскрибированные записи и предоставляет возможность анализировать речевые единицы фонетически (акустически, артикуляционно), просодически и фонологически на уровне звука/фонемы, слова, лексемы, синтагмы, словосочетания и предложения. Существуют три подкорпуса: Wort-Korpus, Fix-Korpus и Frei-Korpus. Wort-Korpus содержит записи отдельно произнесенных слов на разных славянских языков (в настоящее время больше всего на сербском, хорватском и боснийском). Fix-Korpus охватывает записи озвученных текстов. Frei-Korpus состоит из записей спонтанной речи и их транскриптов.

Wort-Korpus

Speech Korpus

The screenshot shows the Wort-Korpus search interface. A search bar contains 'balon'. Below it, a note says 'z.B. balon, ba*, m*'. The search results table has columns: Suchergebnisse, Treffer: 1 - 3 von 3, Wort, m/f, Nationalität, Geburtsjahr, Ort, Geburtsort, Wohnort, Muttersprache. The results are:

	Wort	m/f	Nationalität	Geburtsjahr	Ort	Geburtsort	Wohnort	Muttersprache
1	balon	f	hrvatska	1983	Zenica	Bosna i Hercegovina	Graz	hrvatski
2	balon	f	hrvatska	1980	Rotterdam	Holandija/Nizozemska	Ostjek	hrvatski
3	balon	f	hrvatska	1977	Bjelovar	Hrvatska	Bjelovar	hrvatski

Фонетическая и просодическая транскрипция проходит верификацию экспертов (как минимум одного) при помощи специальной программы –

Valorisarium'a. Speech-Korpus связан с программой Akzentarium, предоставляющей возможность получить информацию об ударении (в данный момент для более чем 120 000 слов сербского, хорватского и боснийского языков) и позволяющей ставить ударение и получать информацию о внесенных в нее словах. В Transkriptarium'e дается информация о том, как произносятся слова, а Suprasemgentarium указывает на интонационную структуру предложения.

Frei-Korpus является коллекцией свободно/спонтанно произнесенных текстов, относящихся ко всем функциональным стилям (литературно-художественному, публицистическому, официально-деловому, публицистическому, научному и разговорному).

Fix-Korpus

Speech Korpus

The screenshot shows the Fix-Korpus search interface. It includes fields for Geschlecht (both), Muttersprache (srpski), Geburtsort (alle), Wohnort (alle), and Chiffre (alle). A search button labeled 'suchen' is at the bottom right.

Suchergebnisse: Jutros su me vrlo rano ptice probudile.

Treffer: 1 - 15 von 30

	m/f	Nationalität	Geburtsjahr	Ort	Geburtsort	Staat	Wohnort	Muttersprache
1	m	srpska	1969	Užice	Srbija	Beograd	srpski	
2	m	srpska	1931	Beograd	Srbija	Beograd	srpski	
3	f	srpska	1974	Niš	Srbija	Graz	srpski	
4	m	srpska	1988	Novi Kneževac	Srbija	Zmajevac	srpski	
5	f	srpska	1987	Zrenjanin	Srbija	Zrenjanin	srpski	
6	m	srpska	1987	Beograd	Srbija	Stara Pazova	srpski	
7	f	srpska	1967	Srpski Itebej	Srbija	Novi Sad	srpski	
8	m	srpska	1938	Virovitid	Bosna i Hercegovina	Sarajevo	srpski	
9	m	srpska	1954	Kalinovik	Bosna i Hercegovina	Istočni Sarajevo	srpski	
10	f	srpska	1923	Virovitid	Bosna i Hercegovina	Istočno Sarajevo	srpski	
11	m	srpska	1931	Dragelji	Bosna i Hercegovina	Sarajevo	srpski	
12	f	crnogorska	1981	Bar	Crna Gora	Bar	srpski	
13	f	crnogorska	1985	Kotor	Crna Gora	Herceg Novi	srpski	
14	m	srpska	1939	Drvra	Bosna i Hercegovina	Beograd	srpski	
15	m	srpska	1950	Kalinovik	Bosna i Hercegovina	Podgorica	srpski	

В Speech-Korpus'e пользователю предоставляется несколько возможностей в поиске: по отмеченным предложениям, по частям с транскрипцией, акцентуацией и супрасегментацией и т. д.

Для обработки данных о записях используются особые процедуры, начинающиеся с заполнения бланка с основными персональными данными. Он разделен на три части: Опрошенный, Запись и Корпус. В первую графу вносится основная информация о говорящем (пол, национальность, религия, год и место рождения, место жительства, профессия, место обучения, родной язык, важнейшие места проживания, иностранный язык и т. п.), дата и место записи, во второй следуют данные о теме, месте, обстоятельствах записи, оборудовании, длине и формате записи, операторе, а в третью вносятся данные о том, для какого типа исследования использу-

ется записанный материал (для устного, спектрального или какого-нибудь другого анализа).²³ В конце бланка находится графа для письменного согласия говорящего на публичное представление (за определенным шифром). В этих целях выбирается пароль, состоящий из названия места записи, трехзначного номера и букв, обозначающий язык (напр., **b**, **k** **s** указывают на боснийский/бошняцкий, хорватский или сербский) типа split_001k. После внесения пароля открывается полный бланк с данными об опрошенном и о записи. Если нужно добавить дополнительную информацию, скажем, о месте или языке, нужно войти в графу „Внести новое место или пароль“.

9. Текстуальный Гралис-Корпус (Gralis Text-Korpus) дает возможность спаривать тексты на всех Сляз, причем в настоящее время основной уклон при его создании делается на пополнении материала по южнославянским языкам (Bg, Bs, Hr, Mk, Mo, Sl, Sr) и крупнейшему славянскому языку (Ru). Благодаря созданной инфраструктуре можно проводить параллелизацию (а) всех славянских языков и немецкого языка, (б) славянских языков по отдельным группам (восточнославянской, западнославянской, южнославянской), (в) языков только из одной группы (например, южнославянской). В рамках Текстуального Гралис-Корпуса существуют индивидуальные подкорпусы. Пока разработаны, четыре из ряда запланированных: корпус Иво Андрича (1892–1975, лауреата Нобелевского премии за литературу), Бранко Чопича (1915–1984, величайшего славянского повествователя, юмориста и сатирика), Зорана Живковича (род. 1948, в настоящее время самого переводимого писателя бывшей Югославии) и Блаце Конеского (1921–1993, виднейшего македонского писателя, поэта и филолога). Особый тип составляет образовательный корпус, в состав которого входят тексты, необходимые дипломантам, аспирантам, докторантам и диссертантам для подготовки их научных работ. Со всеми подкорпусами Текстуальный Гралис-Корпус содержит около шести миллионов словоупотреблений. Для свободного использования он открыт в той части, в которой решен вопрос об авторских правах.

В Text-Korpus включаются тексты всех функциональных стилей и сопровождаются основной метаязыковой и грамматической аннотацией. Gralis-Korpus имеет три типа разметок: метатекстуальную, экстралингвистическую и морфосинтаксическую. Последняя и поиск по ее параметрам проводится при помощи особой программы – „MorphoGenerator“.

²³ В анализе устного материала используются два метода: акустический (при помощи программы типа Praat измеряется расстояние между Сляз по основным параметрам: длине произнесенных звуков, изменению в интонации и т. п.) и слуховой.

most			mosti mostovi			
Kasus	Singular		Plural			
		m			m	
Nominativ	most	N-msn----	21Sm-08	mosti	N-mpn----	21Sm-08
Genitiv	mosta	N-msg----	21Sm-08	mostovi	N-mpg----	21Sm-08
Dativ	mostu	N-msd----	21Sm-08	mostima	N-mdp----	21Sm-08
Akkusativ	most	N-msa--n-	21Sm-08	mostovima	N-mpd----	21Sm-08
Vokativ	moste	N-msv----	21Sm-08	mosti	N-mpv----	21Sm-08
Instrumental	mostom	N-msi----	21Sm-08	mostovi	N-mpv--21Sm-08	
Lokativ	mostu	N-msl----	21Sm-08	mostima	N-mpi----	21Sm-08
				mostovima	N-mpi--21Sm-08	
				mostima	N-mpl----	21Sm-08
				mostovima	N-mpl--21Sm-08	

С текстуальным корпусом связан онлайн-словарь „Lexikarium“ и программа „Akzentarium“, позволяющие получать комплексную просодическую, лексикосемантическую и грамматическую информации.²⁴ Корпус предоставляет два типа поиска: простой и сложный (при помощи т. наз. CQP-синтаксиса). Результаты поиска отмечаются желтым или синим цветом. Над полученными предложениями стоит стрелочка, указывающая на подлинник (если речь идет об интернет-адресе, приводится ссылка). Из корпуса можно извлекать частотные списки.

Мост на Дрине	
Russisch	Там, где Дрина всей тяжестью своей зеленой и вспененной водной лавины извергается как бы из сомнущей стены отвесных черных гор, стоит большой каменный мост строих пропорций с одиннадцатью широкими пролетами
Serbian	На том месту где Дрина из једног целом свеју воду уздече масе , зелена запенјена , из првог затвореног стакла смрт струни планина , струја из једног лукова строгог резана
Deutsch	An dieser Stelle , wo die Drina mit dem ganzen Gewicht ihrer Wassermassen , grün und schläumend , aus dem steinernen geschlossenen Massiv der schwarzen und steilen Berge hervorbricht , steht eine große , gleichmäßig geschnittenne Brücke aus Stein mit elf weitgespannten Bögen
Bulgarsки	На това място , дето Дрина с цялата тежест на своите зелени и пенливи водни маси избива от затворения нагред масив на черните и стръмни планини , има голям , стройно изграден от камък мост с единадесет широки свода
Mazedonisch	На то място каде што Дрина изливнува со целата тежина на своите водни маси , зелена и запенета , од првндно затворениот склон на црните и стръмни планини , стои еден голем , складно издаден мост од камен , со единадесет широки арки .

Создание письменного Гралис-Корпуса преследует цель (а) собрать и обработать материал для Сляз, необходимый для анализа отношений

²⁴ В качестве дополнительного средства для получения и изучения материала можно использовать онлайн-программу по созданию, проведению и обработке опросов под названием „Гралис-Анкетариум“ (разработанный в рамках FWF-проекта P19158-G03), а также программу „Гралис-Прескриптиарий“ по изучению орфографических норм исследуемых языков.

между ними, их совпадений, сходств и различий, а также расстояния. Это подразумевает проведение следующих корпусных процедур: **1.** количественное и качественное сопоставление категорий параллельных фрагментов и единиц текстов, **2.** анализ статистики пар относительно эквивалентных единиц параллельных текстов каждого из языков по ряду параметров (например, по статистике типов соответствия: слово языка¹ – слово языка², слово языка¹ – словосочетание языка², словосочетание языка¹ – слово языка², слово языка¹ – Ø языка² и т. п.), **3.** интерпретация, качественный анализ статистических данных, **4.** выработка и согласование параметров для последующей корректной сопоставимости получаемых результатов (например, составление частотных словарей словоформ, лексем, морфологических категорий и т. п. для каждого из текстов, для текстов определенного жанра, функционального стиля и всего корпуса текстов данного языка в целом), **5.** исследование внутренних характеристик текстов каждого из языков (интракорреляции) и сопоставление корпусов по ним, а также по типам относительно эквивалентных единиц текстов.

Что касается последовательности процедур, материал сначала вносится в „Roh-Korpus“ („Сырой корпус“) на одном из самых больших серверов Университета Граца (Gedra). Тексты, находящиеся в интернете и не нуждающиеся в решении вопроса авторских прав, включаются в „Warte-Korpus“ („Корпус в ожидании“), расчлененный на пять частей для каждого функционального стиля („FS-Korpus“): литературно-художественный, публицистический, научный, официально-деловой и разговорный. В рубрике „Мета-Корпус“ („Meta-Korpus“) рассматриваются теоретические и практические вопросы создания корпуса и проводится дискуссия между сотрудниками проекта.

Обработка материала проходит в два этапа. На первом отдельно готовятся тексты для каждого славянского языка. Важнейшей частью такой работы является аннотация (паспортизация) – метаязыковая и грамматическая. Метаязыковая аннотация состоит из указания сведений об источнике (авторе, заглавии текста, месте и где издания, количестве страниц, издательстве, переводе и др.). Грамматическая аннотация указывает на морфологическую структуру, словоизменение и тип сочетаемости с другими словами в рамках словосочетания и предложения. Существуют различные методы и кодировки, используемые для грамматического аннотирования. Так как ГРАЛИС-Корпус требует полностью унифицированной аннотации (чтобы искать информацию для всех СЛЯЗ), необходимо выбрать ту, которая больше всего годится для изучаемых языков. Опыт и проведенная работа в процессе подготовки параллельного корпуса для Bs, Hr Sr (Gralis BKS-Korpusa) свидетельствуют о том, что в этих целях самым целесообразным является использование кодировки Multext-East (Multilingual

Texts and Corpora for Eastern and Central European Languages – multilingual dataset for language engineering research and development: MultiText East-www), разработанной в 2004 году группой авторов во главе с Томажем Эрявцем. Так как в рамках Multext-East существует унифицированная система грамматического аннотирования для ряда славянских языков (Bg, Cs, Hr, Ru, Sl, Sr), она будет использована и приспособлена для всех других СЛЯЗ. После проведения разметки текст расчленяется на предложения, в результате чего получится система, в которой каждому предложению языка А соответствует предложение языка В, С... Параллелизация состоит в том, что тексты двух или более СЛЯЗ объединяются, потом проводится их выравнивание. Если, например, в одном языке абзац состоит из трех предложений, а в другом из пяти, приходится устраниить такое неравновесие. Для автоматизации данного процесса используются имеющиеся (не)модифицированные разработки. Если существующие программы не в состоянии выполнить поставленные задачи, создаются новые инструментальные средства, позволяющие, в частности, в одном комплексе объединять тексты различных СЛЯЗ, автоматически находить несовпадения в числе предложений и, насколько это возможно, автоматически делать исправления. На этом заканчивается параллелизация и начинается серверная работа, для чего используются IMS Corpus Workbench (CQP) и Asset-Management. Этот этап состоит в том, что из корпусного материала создаются списки языковых единиц, которые затем превращаются при помощи программы MySql в реляционные базы данных, на основании которых можно готовить словари различного типа. IMS Open Corpus Workbench представляет собой набор средств для администрации, подготовки и осуществления поиска в больших текстовых корпусах с лингвистической аннотацией. Его главным компонентом является гибкая и продуктивная поисковая программа CQP (Corpus Query Processor). Первоначально разработанный в Институте машинной обработки языка Штутгартского университета, он был в 2007 году выпущен как программа с открытым кодом (open-source software) с GPL лицензией (GNU General Public License) и размещен на SourceForge. CWB использует для хранения корпуса свой собственный формат. Быстрый доступ достигается за счет бинарной кодировки, полный индекс способствует эффективному поиску словоформ и аннотаций; используются специальные алгоритмы сжатия. В зависимости от аннотации размер корпуса может достигать 500 миллионов словоупотреблений (tokens). CWB содержит инструменты для кодировки, индексации, сжатия, декодирования и частотных распределений, общий реестр, в котором хранится информация о корпусе (название, атрибуты, место нахождения), поисковую программу (CQP), осуществляющую быстрый поиск с использованием синтаксиса регулярных выражений по

значениям атрибутов отдельных позиций (например, по морфологическим тэгам). Поисковая система предоставляет простой и расширенный поиск, который базируется на CQP-синтаксисе (Suche mit CQP-Syntax) и предлагает очень широкие и разнообразные комбинации. Результаты поиска выводятся на монитор вертикально. Название текста/источника обознается желтым цветом. Нажатием стрелочки с левой стороны заглавия можно получить основную метаинформацию (об авторе, месте и времени издания, издательстве, числе страниц и т. п.). При отображении результатов поиска пользователь может определить размер отображения. Допускаются различные виды сортировки строк, подсчитываются частоты (например, для комбинации слов), составляется многоязычный индекс для параллельных корпусов.

Корпусная работа подразумевает: 1. выработку инструментальных средств для создания корпуса, его ведения (пополнения, предобработки текстов, паспортизации, контроля их параметров и т. п.), структурирования (разметки структурными пометами), категоризации (качественной квалификации его фрагментов и единиц, выделенных в ходе структурирования текстов), выравнивания и фиксации связи между относительно эквивалентными элементами параллельных текстов корпуса, 2. создание, тестирование, наполнение, пробную эксплуатацию и развитие и совершенствование оболочки (программного средства) для выполнения работ по созданию, ведению, структурированию, категоризации, выравниванию элементов параллельных текстов, 3. выработку и согласование схемы метаязыка, а также структурной и категориальной разметки фрагментов и единиц текстов, 4. переработку текстов каждого из языков, контроль их параметров и характеристик паспортизации, структурную разметку текстов, категориальную разметку фрагментов и единиц текстов (как минимум планируется осуществить лемматизацию и морфологическая квалификацию словоформ), выравнивание (Alignment) соответствующих эквивалентных фрагментов параллельных текстов, 5. выработку и согласование функционально-стилистической (жанровой) схемы сбора параллельных текстов для корпуса, 6. сбор и исходную паспортизацию (метаразметку) параллельных текстов по каждому из языков.

10. В процессе создания корпуса кроме подготовки базовых текстов и аннотации метаданных в XML-формате необходимо переработать цифровые ресурсы (MP3-, WAV-данные и т. п.) в рамках Speech-Korpus'a. Asset Management System (AMS) дает такую возможность в форме Workflows.²⁵ В узком смысле этим понятием обозначаются системы хранения, управления и подготовки цифровых ресурсов, накапливаемых в большом

²⁵ Развитие устного корпуса (Speech-Korpus) требует больших усилий, много времени и значительных финансовых средств.

количестве. В отличие от Content Management System, в Asset Management выступает на передний план, прежде всего, архивирование постоянных метаданных, доступных для цитирования и гибкого использования цифровых ресурсов. Основная идея Asset Management состоит в том, чтобы раздробленные цифровые ресурсы получили одну центральную ИТ-структуру и этим обеспечили продолжительное архивирование имеющегося цифрового фонда знаний, доступного для цитирования. С (текстуально-)технологической точки зрения Asset Management System предоставляет, в первую очередь, очень гибкий способ хранения документов, основанных на XML, и управление корпусом. Именно из-за этих свойств AM-System являются надежным местом для хранения корпусных текстов.

В качестве платформы для трансфера материала используется Open Source Projekt „Fedora“ (Flexible Extensible Digital Object Repository Architecture), разработанный в Cornell University (который с 1997 г. вместе финансировали Виргинский университет и Mellon Foundation, а ранее – National Science Foundation). „Fedora“ предоставляет структуру для сохранения цифровых web-ресурсов и управления ими (Repository). Она основана на архитектуре SOA. Речь идет о распределющейся системной архитектуре, не зависящей от платформы и основанной на web-сервисе (SOAP, Simple Object Access Protocol). Она содержит полный текстовый регистр (VolltextIndex), базируемый на Apache Lucene, и версионное управление (Versionsverwaltung) Asset-содержанием (Assetinhalte). Архитектура SOA обладает регистром метаданных с языком запросов ITQL типа SQL (Tucana Technologies), основанном на RDF, и четкой адресацией цифровых ресурсов на базе URL. Для нее характерны тонко расстраиваемые права доступа к Assets и к его частям на основании XACML (Extensible Access Control Markup Language), а также стандартные форматы импорта и экспорта: METS (Metadata Encoding and Transmission Standard). Существует поддержка стандартных протоколов обмена метаданными типа OAI-PMH. На основе носящей технологии (Trägertechnologie) Apache Tomcat можно также осуществить при помощи Repository Clustering загрузочную балансировку (Load Balancing) системного окружения с соответствующим числом пользователей (Concurrent Users).

Для успешного и эффективного проведения корпусных работ существует особый on-line Projekt-Management, в рамках которого созданы базы данных „Персоналиум“, „Библиотекариум“, а также „Слав-Форум“. „Персоналиум“ содержит все основные данные о сотрудниках проекта. В эту базу данных вносятся тексты для публикаций, резюме, хендауты и презентации для совместных конференций. В „Библиотекариуме“ стоит в распоряжении литература, относящаяся к исследуемой теме. Для постоян-

ного контакта и быстрого обсуждения определенных вопросов используется Slawistik-Forum (<http://www-gewi.uni-graz.at/gralis/slau/forum>).

11. В работе над Гралис-Корпусом ожидаются следующие результаты: 1) создание параллельного устного и письменного корпусов для всех Сляз, снабженного метазыковой и грамматической аннотацией (объем корпуса, широта охвата языков и глубина аннотированния будут зависеть от уровня финансирования), охватывающего все функциональные стили и доступного для всех или как можно большого числа пользователей, при полном соблюдении авторских прав, 2) создание реляционных баз данных – „Slaw-Lexikarium“ с корпусными словарными единицами, „Slaw-Grammatikarium“ с корпусными грамматическими формами, „Slaw-Akkzentarium“ с уда-рениями для всех слов в корпусе, „Slaw-Präskriptarium“ с правилами для написания корпусных слов и „Slaw-Bibliothekarium“, охватывающий важнейшую литературу.

Литература

- Ammon, U. 1986. Explikation der Begriffe 'Standardvarietät' und 'Standardsprache' auf normtheoretische Grundlage. In: Holtus G., Radtke E. (Hg.) *Sprachlicher Substandard*. Tübingen: Max Niemeyer Verlag, 1–63.
- Ammon, Ulrich. 1987. Funktionale Typen/Statustypen von Sprachsystemen. In: Ammon U., Dittmar N., Mattheier K.J. (Hg.) *Sociolinguistics/Soziolinguistik*. First Volume / Erster Halbband. Berlin, New York: Walter de Gruyter, 230–263.
- Hudson, R. A. 1990 [1980]. *Sociolinguistics*. Cambridge, New York, Portchester, Melbourne, Sydney: Cambridge University Press.
- Tošović Branko 2001: *Korelaciona sintaksa*. Projektional. Graz: Institut für Slawistik der Karl-Franzens-Universität.
- Tošović B. 2002. *Funkcionalni stilovi. Funktionale Stile*. Graz: Institut für Slawistik der Universität Graz.
- Tošović, B. 2008. *Die Unterschiede zwischen dem Bosnischen/Bosniakischen, Kroatischen und Serbischen*. Wien, Berlin: LIT.
- Tošović, B. (Hg). 2008a. Das Gralis-Korpus. In: *Die Unterschiede zwischen dem Bosnischen/Bosniakischen, Kroatischen und Serbischen*. Münster: LIT, 724–827.
- Tošović, B. 2008b. Das Gralis-Akkzentarium. In: *Die Unterschiede zwischen dem Bosnischen/Bosniakischen, Kroatischen und Serbischen*. Münster: LIT, 770–776.
- Tošović, B. 2008c. Das Gralis-Bibliothekarium. In: *Die Unterschiede zwischen dem Bosnischen/Bosniakischen, Kroatischen und Serbischen*. Wien, Berlin: LIT, 807–812.

- Tošović, B. 2008d. Das Gralis-Präskriptarium. In: *Die Unterschiede zwischen dem Bosnischen/Bosniakischen, Kroatischen und Serbischen*. Wien, Berlin: LIT, 822–825.
- Tošović, B. 2008e. Der Unterschied. In: *Die Unterschiede zwischen dem Bosnischen/Bosniakischen, Kroatischen und Serbischen*. Wien, Berlin: LIT, 143–185.
- Tošović B. 2008f. Сопоставительное изучение славянских языков при помощи многоязычного „Гралис-Корпуса“. In: *Izučavanje slovenskih jezika, književnosti i kultura kao inoslovenskih i stranih*. Beograd: Slavističko društvo Srbije, 336–340.

Электронные корпусы

- Chemnitz German-English Translation Corpus:
<http://www.tu-chemnitz.de/phil/InternetGrammar>.
- Compara-www: <http://www.linguateca.pt/COMPARA/index.php>.
- EPC-www: <http://www.statmt.org/europarl>.
- Gralis-Korpus-www:
<http://www-gewi.kfunigraz.ac.at/gralis/0.Projektarium/Gralis-Korpus/korpus.html>.
- Infostream-www: <http://ling.infostream.ua/>.
- Kacenka: <http://www.phil.muni.cz/angl/kacenka/kachna.html>.
- Kollokation-www:
<http://www.kokken.go.jp/public/world/mirror/www.ids-mannheim.de/gra/kollokation.html>.
- Leeds-www: <http://corpus.leeds.ac.uk>.
- Lilabar-www: <http://lilabar.com/index.php>.
- Maastr-www: <http://www.philhist.uni-augsburg.de/lehrstuehle/anglistik/>.
- Opus-www: <http://opus.lingfil.uu.se>.
- RNK-www: <http://corpora.yandex.ru/search-para.html>.
- RPC-www: <http://www.uni-regensburg.de/sprache-literatur-kultur/slavistik/rundums-institut/regensburger-korpora/index.html>.
- Ruscorpora-www: <http://ruscorpora.ru>; <http://ruscorpora.ru/search-para.html>.
- SPI-www: <http://nevmenandr.net/slovo>.
- Traumdeutung-www: http://www.aac.ac.at/lab_parallel_freud.html.