

# Использование лексико-грамматических баз данных в русской диалектной лексикографии<sup>1</sup>

## The use of lexico-grammatical databases in the Russian dialectal lexicography

**Тер-Аванесова А. В.** (teravan@mail.ru)

Институт русского языка им. В. В. Виноградова РАН

**Крылов С. А.** (krylov-58@mail.ru)

Институт востоковедения РАН, Институт системного анализа РАН

С помощью СУБД STARLING обогащена построенная ранее лексико-грамматическая база данных (ЛГБД) по русским народным говорам с различением двух фонем «типа о». К созданной ранее базе данных по среднерусскому говору с. Пустоша Шатурского р-на Московской обл. добавлена ЛГБД по вологодскому слободскому говору, включающая ок. 30 тыс. словоформ, представляющих ок. 4500 лексем. Ядерный диалектный корпус (ЯДК) содержит тексты с частичной лексико-грамматической разметкой. В сентенциальной базе единицами являются предложения ЯДК в фонологической транскрипции, пронумерованные в порядке вхождения в ЯДК. На ее основе создан прямой алфавитный лексико-грамматический конкорданс и обратный алфавитный лексико-грамматический указатель словоформ. ЛГБД содержит информацию об условной фонологической транскрипции данной единицы, о словоизменительных и акцентных типах лексем, смысловые пометы о лексических значениях семантических диалектизмов, а также метаязыковые социолингвистические пометы о возрастных и территориальных особенностях употребления словоформы.

### 1. Предмет исследования и материал: русские народные говоры с различением двух фонем «типа о»

В рамках проекта РГНФ в 2006 г. было продолжено создание ЛГБД по русским народным говорам с различением двух фонем «типа о». К созданной ранее базе данных среднерусского (владимирско-

поволжского) говора с. Пустоша Шатурского р-на Московской обл. добавилась построенная в формате STARLING лексико-грамматическая база данных по севернорусскому слободскому говору деревень Арзубиха, Захариха и Злобиха Харовского р-на Вологодской области.

Предметом исследования являются русские говоры, в системе вокализма которых представлены две фонемы «типа о», распределенные в соответствии с правилом Л.Л. Васильева — А.А. Шахматова: «о закрытое» (фонема /о/) выступает на месте \*о под праславянским «восходящим» ударением, «о открытое» (фонема /о/) — на месте \*о под «нисходящим» ударением, на месте \*ъ, \*е, \*ь. В говорах с различением двух фонем «типа о» обычно также различаются две фонемы «типа е», наряду с фонемами /а/, /у/, /и/, в связи с чем их системы вокализма получили название семифонемных. В настоящее время такие говоры достаточно редки, не образуют сплошных ареалов, сохраняются главным образом в восточной части Европейской территории России и лишь отдельными вкраплениями — к югу и юго-западу от Москвы. Данные русских говоров с семифонемным вокализмом имеют особое значение для истории русского языка и славянской акцентологии, поскольку тембр ударного о < \*о является (по крайней мере, в части случаев) отражением праславянских слоговых тонов.

Некоторые косвенные данные свидетельствуют о том, что в прошлом системы вокализма рассматриваемого типа были распространены в русских говорах гораздо шире. Ареалы таких систем должны были быть не меньше современных ареалов нескольких типов диссимилятивного яканья, пред-

<sup>1</sup> Данная работа выполнена при финансовой поддержке Программы ОИФН РАН «Генезис и взаимодействие социальных, культурных и языковых общностей» (проекты «Восточнославянский диалектный корпус: праславянское наследие и лингвогеография» и «Генезис балто-славянской языковой общности: акцентологический аспект»), а также гранта РГНФ № 08-04-12132в.

полагающих семифонемный и шестифонемный ударный вокализм (обоянский, задонский, дмитриевский, новосёлковские, ореховские типы яканья). Карты и материалы Диалектологического атласа русского языка показывают, что небольшие кружевные ареалы юго-восточных семифонемных систем вокализма «вписаны» в несравненно более обширные ареалы перечисленных типов диссимильятивного яканья. Следовательно, эти типы яканья являются косвенным свидетельством наличия в прошлом различия под ударением двух фонем «типа о» на гораздо большей территории, чем сегодня. Локализация памятников письменности XIV–XVII вв., графико-орфографические системы которых отражают противопоставление двух фонем «типа о», показывает, что говоры с различием двух о в старорусский период были распространены почти на всей территории русского языка.

В 2006 г. были проведены три диалектологические экспедиции для сбора материала по говорам изучаемого типа: в с. Новосёлки Рыбновского р-на Рязанской обл., с. Пустоша Шатурского р-на Московской обл. и в дд. Арзубиха, Захариха и Злобиха Харовского р-на Вологодской обл. Расшифровки магнитофонных записей речи уроженцев названных вологодских деревень стали материалом для ЛГБД.

Несколько говором рассматриваемого типа были предварительно сопоставлены, в том числе с помощью построенных ЛГБД, как в отношении акцентных систем и распределения двух фонем «типа о», так и в отношении их лексического состава: западно-вологодские слободской и тотемский (последний — по описанию О. Брока); владимирско-поволжский говор с. Пустоша; восточный средне-русский акающий говор с. Лека Шатурского р-на Московской обл. (по описанию А.А. Шахматова); задонский говор (по материалам В. Тростянского); рязанский говора с. Новосёлки.

Говоры с противопоставлением двух о обнаруживают сильные различия по общим наборам признаков (они относятся к разным наречиям и группам говором). Одновременно, относясь к восточной диалектной зоне, все говоры с двумя о имеют ряд важных общих черт: «моновариантное» склонение типа рус. лит., маргинальную подвижность ударения в прош. времени глаголов с корнями на нешумные, в целом схожее распределение непроизводных существительных по акцентным типам и ряд других сходств. Так, все семифонемные говоры обнаруживают нетривиальное сходство: сохранение у небольшого числа существительных муж. рода (\*и-, \*i- и консонантные основы а. п. d) рефлекса смешанной акцентной парадигмы (с формой-энклитоменом в И.ед. и окситонезой прочих форм). Обнаружены признаки, противопоставляющие друг другу отдельные группы говором с различием двух фонем о, например, 1) /yo/ из \*o в формах мн. числа

слов ж. и ср. рода а. п. b (вдубы, вдубами; долубы, долубами) в средне- и южнорусских говорах рассматриваемого типа; в севернорусских в тех же случаях — /o/; 2) накоренное ударение в наст. времени и пов. наклонении i-глаголов а. п. b при насуффиксальном — в инфинитиве и прош. времени, характерное для говоров Рязанской группы и «рязанского ареала» говором с различием двух о (Пустоша хобжу, худдишь, Новосёлки хубжу, худдишь). Списки глаголов с указанной инновацией, однако, сильно различаются в говорах Пустошей и Новосёлок: если в Пустошах этот список ограничен итеративами а. п. b1 (ходить, носить, возить, водить, молотить, просить и т. д.), то в Новосёлках в него входят каузативы и деноминативы а. п. b2 и даже с. Последнее различие должно указывать на гетерогенный характер говором с различием двух фонем «типа о» в «рязанском ареале».

Построение лексико-грамматической базы данных слободского говора (Харовский р-н Вологодской обл.). Аудиозаписи речи носителей говора старшего поколения были расшифрованы и записаны в аллофонемной транскрипции. На основе получившихся текстов с помощью интегрированной информационной среды STARLING (автор — чл.-корр. РАН С. А. Старостин) построена лексико-грамматическая база данных говора — так называемый ядерный диалектный корпус (ЯДК). ЯДК представляет собой исчерпывающее описание говора в рамках определенного корпуса текстов и охватывает тексты общей длиной около 30 тыс. речевых словоформ. Они репрезентируют 7047 языковых словоформ без учёта пунктуации, 9591 пунктуационно-грамматическую словоформу (пунктуационный вариант языковой словоформы).

## 2. Структура базы данных слободского говора

Структура базы данных слободского говора идентична структуре созданной ранее базы данных говора с. Пустоша Шатурского р-на Московской обл. В качестве исходной базы данных выступает ЯДК. Лингвистическая информация в ЯДК организована по многоступенчатому принципу. Выделяется 7 уровней членения письменного текста; на каждом из них выделяется своя основная (базовая) единица членения; каждой единице членения каждого уровня в ЯДК приписан уникальный номер, способный служить адресом отсылки к этой единице.

1. Уровень целого текста. На этом уровне вводятся параметры, характеризующие личность информанта: фамилия, имя, отчество, год и место рождения, образование и т.п.

2. Уровень абзаца (сверхфразового единства). Сверхфразовое единство — это отрезок текста, пун-

ктуационно выделенный особым абзацным делимитатором («красной строкой», «отступом»). У сверхфразового единства есть некоторая единая общая смысловая тема.

3. Уровень предложения (сентенциальный уровень). Границы предложений помечены сентенциальными делимитаторами. В начале предложения стоит инициальный делимитатор — суперсегментная пунктограмма «заглавности»; в конце предложения стоит финальный делимитатор — пунктограммы «.», «?», «!», «...». Содержательно предложение соответствует законченной мысли, а фонетически — интонационно законченному отрезку.

4. Уровень клаузы // предикации (клаузальный уровень). Границы клауз помечались так: предложение состоит из клауз, а между клаузами внутри предложения стоит один из клаузальных делимитаторов. К ним относятся пунктограммы «;», «:» и «—». Содержательно и интонационно клаузы примерно соответствуют простым предложениям и отдельным предикациям в составе сложных предложений.

5. Уровень синтагмы (синтагматический уровень). Границы синтагм внутри клаузы помечены пунктуационным синтагматическим делимитатором — пунктограммой «запятая». Содержательно и интонационно синтагмы примерно соответствуют словосочетаниям.

6. Уровень такта (тактовый уровень). Границы такта в ходе расшифровки помечались стандартной орфографической пунктограммой «пробел», но после создания ЯДК они были размечены вручную так: был использован пунктуационный тактовый делимитатор — пунктограмма «знаменательный (паузальный) пробел». Такты примерно соответствуют фонетическим словам, членам предложения, «синтаксическим молекулам», формам слова (как аналитическим, так и синтетическим). Важнейшее фонетическое свойство такта: внутри него невозможна (или по меньшей мере нетипична) пауза.

7. Уровень глоссы (глоссовый уровень). Границы глоссов в ходе расшифровки помечались либо стандартной орфографической пунктограммой «пробел», либо стандартной орфографической пунктограммой «дефис», но после создания ЯДК их границы были размечены вручную. Каждый такт состоит из одного или нескольких глоссов. Глоссы, входящие в состав одного такта, обладают признаком потенциальной подвижности в предложении. Для обозначения границ глоссов при разметке был использован специальный набор нескольких метаязыковых глоссовых делимитаторов — пунктограммы «служебных пробелов». Выделены служебные пробелы шести типов: «{» между проклитикой и её правой опорой; «}» между энклитикой и её левой опорой; «<» между проклитикоидом и его правой опорой; «>» между энклитикоидом и его левой опорой; «<>» между членами квази-композиата с неустойчивым просодическим центром; «&» между компонентами

«фразеологического штампа» с множеством просодических центров. Глоссы примерно соответствуют по длине морфологическим словам (в т. ч. служебным словам, синтетическим формам слов и подвижным компонентам аналитических форм). Внутри глосса (так же как внутри такта) невозможна пауза. Фактически наиболее близкий аналог глоссов в русском письменном тексте, записанном по правилам русской орфографии — это графические слова.

Ценность предложенной многоуровневой схемы ЯДК состоит в том, что при необходимости вывести на обзор список отрезков текста, обладающих некоторым общим свойством, STARLING позволит пользователю по выбору вывести (на экран, на принтер или в файл) отрезок не только одного формата, но разных форматов — графическую словоформу (глосс), минимальный контекст этой словоформы (аналитическую форму, например, предложно-падежную, сочетание клитики с акцентно автономной словоформой и т. п. — такт), словосочетание (синтагму), предикацию (клаузу), предложение, абзац.

Лингвистическая информация о единицах текста на данном этапе в ЯДК такова: 1. Условная фонологическая транскрипция данной единицы (в сочетании с её пунктуационной разметкой). 2. Словоизменительный и акцентный тип данной единицы. 3. Смысловые пометы (при лексических диалектизмах). 4. Метаязыковые социолингвистические пометы о возрастных и территориальных особенностях употребления словоформы.

### 3. Приложение. Образцы словарных статей (иллюстрирующих семантическое поле «позвоночник» в харовском говоре)

**Лён** 1 'шейный отдел позвоночника' <а.т. В>.

*Вот <nom. sg.> лён это самой у ч'елови́ека. Голова́ с позвоно́чником связана, между́ ним <nom. sg.> лён. А зди́ес го́рло. А шéйя это всё вме́сте шéйя и йёс. А у шш́уки-те ни́ету <gen. sg.> лну́-то ётово. Какóй у шш́уки <nom. sg.> лён. Ни́ету <gen. sg.> лну́ у шш́уки. <Егоров Виктор Никол. 1940 г. р. Род. в д. Злобиха Харовск. р-н, Волог. обл. 7 кл. Зап. Белова, Тер-Аванесова в д. Злобиха, Харовск. р-н, Волог. обл., 2003 г.>.*

**Осёл** 1 'шестой шейный позвонок' <а.т. А>.

*<nom.-acc. sg.> осéw, <gen. sg.> осéла <Клешнина Нина Васил. 1936 г. р. Род. в д. Арзубиха Харовск. р-н, Волог. обл. 7 кл. Зап. Тер-Аванесова в д. Арзубиха, Харовск. р-н, Волог. обл., 2002 г.>.*

*<nom. sg.> Осéw это хря́иш о́коло хрептá, о́коло позвоно́чника, пёрва-та шы́шка. Еши́б до позвоно́чника не дошлó, и шéйя конц'яец'це —*

*ето* <nom. sg.> *осѣѡ*. <nom. sg.> *Осѣѡ боли́т, горі́т, как ц'іре́й рвѣ́т, как пересі́лиш себѣ́, ето́т* <nom. sg.> *осѣѡ*. <Егорова (Фокина) Зинаида Никол. 1941 г. р. Род. в д. Полутиха Харовск. р-н, Волог. обл. 8 кл. Зап. Белова, Тер-Аванесова в д. Злобиха, Харовск. р-н, Волог. обл., 2003 г.>.

**Хребѣ́т** 'хребет, позвоночник без шейного отдела' <а.т. В>.

<nom.-acc. sg.> *хребѣ́т*, <gen. sg.> *хрепта́*, <instr. sg.> *хрептѡ́м*, <loc. sg.> *на хрептіе́*, <nom. pl.> *хрепты́* <Клешнина Нина Васил. ...>; <gen. sg.> *о́коло хрепта́* <Егорова (Фокина) Зинаида Никол. ...>.

**Хрип 1** 'соединение позвоночника и черепа у рыб' <а.т. А sg.>.

*Йа сломі́ѡ шіу́ки* <acc. sg.> *хрї́п, ет то́лько схрў́пало. Фсѣ́, она́ уш готѡ́ѡ. Ф сіе́тку попаде́т, нука́к йейѡ́ задаѡі́т? го́лову рас — фсѣ́. Вѣ́том міе́сте у нейѡ́ сла́бойѡ́ е́то міе́сто-то. А бо́лишў́у ка́к, ника́к немо́жно, мніе́ про́шлој гѡ́т попа́ла, повезло́ — йедва́ <acc. sg.> *хрї́п сломі́ѡ. Шіу́чина бо́лшї́ѡ, ак <acc. sg.> *хрї́п йедва́ сломі́л. Ни́е́тулнў́ у шіу́ки. Йе́сли бы бїѡ́ у шіу́ки лё́н, шіу́ки та́з бы <acc. sg.> *хрї́п не сломі́т. Уш <gen. sg.> *хрї́ па не сломі́т.* <Егоров Виктор Никол. 1940 г. р. Род. в д. Злобиха Харовск. р-н, Волог. обл. 7 кл. Зап. Белова, Тер-Аванесова в д. Злобиха, Харовск. р-н, Волог. обл., 2003 г.>.****

#### 4. Приложение. Полные синонимы в слободском говоре (на материале существительных)

При помощи ЛГБД легко показать, что говор деревень Арзубиха, Захариха и Злобиха, а также других деревень бывш. Слободского с/с Харовского р-на Вологодской обл., является единым. Наблюдаемые различия отчасти объясняются, по свидетельству информантов, как относящиеся к «младшей» или «старшей» разновидностям говора, причем соответствующие единицы «младшей» разновидности, как правило, заимствованы из литературного языка. Имеются и такие различия внутри говора, которые представляют собой внутрисистемные колебания. Крайне редки различия, которые могут претендовать на принадлежность к разным диалектным системам (например, название шеста, вокруг которого укладывают сено в стог: Арзубиха, Захариха *стожаир* или *стогаир*, Злобиха *островишна*; название помещения над избой, чердака: Арзубиха *иизбиця*, Митиха *потолвика*).

Ниже приводятся пары существительных — полных синонимов, выявленных в словаре слобод-

ского говора при помощи ЛГБД. Первый из пары синонимов является элементом традиционного лексического состава говора (иногда даже это — устаревшее слово); второй, как правило, представляет собой заимствование из литературного языка. Этот список выделен из словаря существительных, включающего около 2000 лексем; тем самым, примерно десятая часть словаря существительных говора представляет собой пары полных синонимов.

Пары полных синонимов (в большинстве своих случаев обязанных факту диалектно-литературного двуязычия) выделялись на основе явных показаний информантов («можно так сказать, а можно и так сказать»). Хронологические различия между членами пар («старое» / «новое») отмечались также на основании показаний информантов («сейчас говорят так-то, а раньше говорили так-то»).

*ба́ба — жѡна́, ба́ба — же́ншына, батѡ́к — па́вка, берѣ́мѣ — оха́пка, блюдо́ — ми́ска, бог — ико́на, божѡ́т — хрѣ́сној, божѡ́тка — хрѣ́снама́т, бру́снїца — бру́снїка, брю́шына, брю́хо — жы́вѡт, вар — пи́ена, ве́ред — нары́ѡ, ви́ця — ви́етка, ву́бїха ~ ѡ́ха 'ольха', во́нненця — вѡнну́шка, во́тен — лентя́й, гру́да — ку́чя (предметов), губа́ — подбѡру́ѡдок, губа́ — чя́га, губа́ — пога́нка, дво́йник — близне́ць, долѡ́н — ток (вгумне), ка́таник — ва́ленок, ко́лоба́шка — лепѣ́шка, ко́луѡда — ко́рыто, ком — ломѡ́т (хлеба), ко́сиця — висѡ́к, ко́шуля — шу́ба 'шуба, крытая сукном', ку́т — ку́хня, ку́фты́р — жы́вѡт, избá — кубѡ́ната, и́збиця — черда́к, йи́ежа — йе́да, за́города — и́згород, не́погодь — бу́ря, лабаза́ — лѡса́ (строительные), ла́ѡа — мо́стик, ла́ѡка — магази́н, ле́жен — лентя́й, лё́н — ше́ѡ, ло́шат — ку́ѡн, ля́га — топ, ле́ѡина — болѡ́то, ме́тла — ви́ник, ми́еѡ — за́кѡѡска (из пивного сула), мизгї́р, мызгї́р — па́ѡк, мост — поѡ, мостї́на — полови́ця, набѣ́ре́г — за́кѡѡска (из остатков ржаного теста), на́зе́м — наѡѡ́з, на́стаѡннїця устар. — учї́телнїця, оболѡ́чїна — ту́чя, оболѡ́чїнка — ѡ́блако, обѡ́тка — ѡ́буѡ, остре́жнїк — стрѣ́ха, отѣ́ре́бок — заму́хрышка, отченáш — мо́лїтва, плѣ́нка — пузы́р 'околоплодный пузырь', по́скуѡ́тина, по́скуѡ́тка — па́збишишо, погане́ць — пога́нка, постѣ́лѡя — послїе́д, потѡлу́ѡка — черда́к, простѡкї́ша — простѡкѡ́ѡша, рї́зен — кусо́к (хлеба), роднї́к — ко́луѡдець, ры́ло — ну́бсик (чайника), слї́зен — ули́тка, соба́ка — пѣ́с, стекля́шка, скля́нка — ба́нка, буты́ѡвка, соло́ды́ха — сыройїе́шка, соромѡ́та — сты́д, сіе́ра — смо́ла, ста́ѡа — хлеѡ, тоска́ — бол 'боль', у́лик — у́лей, фата́ — платѡ́к, хребѣ́т — позво́нуѡ́чник (позвоночник без шейного отдела, ср. лё́н), хрѡ́ток — хря́шшык, черѣ́д — ѡ́черед, че́рен — ру́чка (лопаты, вил), чилї́к — побѣ́реѡѡѡтик, чиря́к — чї́рей, ша́м — му́сор, со́р, шы́мора — про́хїнде́й, шубнѡ́к — шу́ба 'шуба мехом вверх'.*

## Литература

1. Брок 1907 — О. Брок. Описание одного говора из юго-западной части Тотемского уезда // Сборник ОРЯС, 1907. — Т. 83.
2. Васильев 1929 — Л. Л. Васильев. О значении каморы в некоторых древнерусских памятниках XVI–XVII вв. К вопросу о произношении звука о в великорусском наречии. Л., 1929.
3. Зализняк 1985 — А. А. Зализняк. От праславянской акцентуации к русской. М., 1985.
4. Крылов С. А. Измерение частотности синтаксических молекул (на материале Генерального корпуса русского языка) // Кибрик А. Е. (ред.), Компьютерная лингвистика и интеллектуальные технологии. Вып. 7 (14). По материалам международной конференции «Диалог'2008» (Бекасово, 4–8 июня 2008 г.), М.: РГГУ, 2008а. С. 254–261.
5. Крылов С. А. О частотном словаре фонетических слов (на материале Генерального корпуса русского языка) // Архипов А. В. и др. (ред.). Фонетика и нефонетика. К 70-летию Сандро В. Кодзасова. М.: Языки славянских культур, 2008б, с. 387–399.
6. Крылов С. А., Тер-Аванесова А. В. Лексико-грамматические базы данных как инструмент диалектологического описания // Труды международной конференции «Диалог 2006». М., 2006. С. 493–497.
7. Ter-Avanesova A. Russian dialects with the distinction of two o-phonemes and their contribution to Slavonic accentology (Русские говоры с различием двух о-фонем и их значение для славянской акцентологии) // Second International Workshop on Balto-Slavic Accentology. Copenhagen, 2006. P. 20–24.
8. Тер-Аванесова А. В. Акцентуационные особенности русских говоров с различием двух фонем «типа о» // Тезисы докладов Международной конференции «Актуальные проблемы русской диалектологии» 23–25 октября 2006 г. М., 2006. С. 177–180.