## Компьютерная лингвистика и интеллектуальные технологии

По материалам ежегодной Международной конференции «Диалог» (2013)

Выпуск 12

В двух томах

Том 1. Основная программа конференции

# Computational Linguistics and Intellectual Technologies

Papers from the Annual International Conference "Dialogue" (2013)

Issue 12

Volume 1 of 2. Main conference program

УДК 80/81; 004 ББК 81.1 К63

# Программный комитет конференции выражает искреннюю благодарность Российскому фонду фундаментальных исследований за финансовую поддержку, грант № 13-06-06047

Редакционная В. П. Селегей (главный редактор),

коллегия: В. И. Беликов, И. М. Богуславский, Б. В. Добров,

Д. О. Добровольский, Л. М. Захаров, Л. Л. Иомдин, И. М. Кобозева, Е. Б. Козеренко, М. А. Кронгауз, Н. И. Лауфер, Н. В. Лукашевич, Д. Маккарти, Й. Нивре, Г. С. Осипов, В. Раскин, И. В. Сегалович,

Э. Хови, С. А. Шаров

Компьютерная лингвистика и интеллектуальные технологии: По материалам ежегодной Международной конференции «Диалог» (Бекасово, 29 мая — 2 июня 2013 г.). Вып. 12 (19): В 2 т.

Т. 1: Основная программа конференции. — М.: Изд-во РГГУ, 2013.

Сборник включает 84 доклада международной конференции по компьютерной лингвистике и интеллектуальным технологиям «Диалог 2013», представляющих широкий спектр теоретических и прикладных исследований в области описания естественного языка, моделирования языковых процессов, создания практически применимых компьютерных лингвистических технологий.

Для специалистов в области теоретической и прикладной лингвистики и интеллектуальных технологий.

© Редколлегия сборника «Компьютерная лингвистика и интеллектуальные технологии» (составитель), 2013

#### Предисловие

12-й выпуск ежегодника «Компьютерная лингвистика и интеллектуальные технологии» содержит материалы 19-й Международной конференции «Диалог». В результате работы 54 рецензентов для сборника было отобрано 84 доклада, охватывающих различные направления исследований в области компьютерного моделирования и анализа естественного языка. В настоящем сборнике представлены:

- Лингвистическая семантика и семантический анализ;
- Формальные модели языка и их применение;
- Теоретическая и компьютерная лексикография;
- Методы оценки (evaluation) систем анализа текстов и машинного перевода;
- Корпусная лингвистика. Создание, применение, оценка корпусов;
- Новые лингвистические ресурсы;
- Интернет как лингвистический ресурс. Лингвистические технологии в Интернете;
- Онтологии. Извлечение знаний из текстов;
- Компьютерный анализ документов: реферирование, классификация, поиск;
- Автоматический анализ тональности текстов;
- Машинный перевод;
- Модели общения. Коммуникация, диалог и речевой акт;
- Анализ и синтез речи.

«Диалог» является ведущей российской конференцией по компьютерной лингвистике и, видимо, единственным в мире форумом, посвященным прежде всего проблемам компьютерного анализа русского языка. Принципиальной особенностью конференции, ее основополагающей традицией является особое внимание к технологиям автоматического анализа текста, основанным на лингвистических моделях. Именно этим объясняется и состав участников, и программа конференции, в которой соседствуют теоретические и прикладные исследования. В «Диалоге» представлены также и работы, сделанные в рамках статистических подходов, что позволяет, в частности, сравнивать полученные результаты.

«Диалог» является не только местом обмена опытом и представления новых достижений. Он является также и форумом для разработки и апробирования методик верификации и оценки как результатов лингвистических исследований, так и эффективности работы различных видов систем анализа текстов на русском языке. Целью этой работы являются единые для авторов и рецензентов принципы доказательства и оценки объективности, эффективности и научной новизны предлагаемых решений и методики проведения сравнительного тестирования, на которых могли бы основываться такие оценки.

Схожие проблемы решает в области информационного поиска семинар РОМИП: не случайно, что вот уже второй год «Диалог» и РОМИП проводят совместные дорожки тестирования, результаты участников которых докладываются на «Диалоге» и публикуются в этом сборнике.

В этом году проводилось два соревнования: по анализу тональности (продолжение тестирования 2012 года) и по оценке систем Машинного Перевода (для англо-русской языковой пары).

Особая роль русского языка обусловливает наличие в программе работ по адаптации к нему известных алгоритмов и методов, разработанных для других языков. Доказанные положительные или отрицательные результаты такого применения рассматриваются рецензентами как новые.

За год, прошедший после последней конференции, «Диалог» понес невосполнимую потерю: ушел из жизни выдающийся лингвист и один из отцов-основателей «Диалога» Александр Евгеньевич Кибрик. Трудно переоценить его роль в создании особой концепции и самой атмосферы конференции, которая сохраняется вот уже почти 40 лет, начиная с первых семинаров середины 70-х годов, из которых и вырос «Диалог». Основными чертами этой концепции всегда оставались широта взгляда, междисциплинарность, сочетание конструктивности и теоретической значимости обсуждаемых проблем. В этом году А. Е. Кибрику посвящено специальное заседание, материалы которого также вошли в сборник.

Несмотря на традиционную широту тематики докладов одного года, они не могут дать полной картины направлений «Диалога». Ее можно получить с помощью сайта конференции www.dialog-21.ru, на котором представлены обширные электронные архивы «Диалогов» последних лет.

Программный комитет «Диалога» Редколлегия ежегодника «Компьютерная лингвистика и интеллектуальные технологии»

#### Организаторы

Буате Кристиан

Раскин Виктор

Хови Эдуард

Сегалович Илья Валентинович

Селегей Владимир Павлович

Шаров Сергей Александрович

Ежегодная конференция «Диалог» проводится под патронажем Российского Фонда Фундаментальных Исследований при организационной поддержке компании АВВҮҮ.

Учредителями конференции являются:

- Институт лингвистики РГГУ
- Институт проблем информатики РАН
- Институт проблем передачи информации РАН
- Компания АВВҮҮ
- Компания Yandex
- Филологический факультет МГУ

Конференция проводится при поддержке Российской ассоциации искусственного интеллекта.

Гренобльский университет

#### Международный программный комитет

руате кристиан	т реноольский университет
Богуславский Игорь Михайлович	Политехнический университет Мадрида
Гельбух Александр Феликсович	Национальный политехнический
	институт, Мехико
Иомдин Леонид Лейбович	Институт проблем передачи
	информации PAH
Кобозева Ирина Михайловна	Филологический факультет МГУ
Козеренко Елена Борисовна	Институт проблем информатики РАН
Корбетт Гревил	University of Surrey, UK
Кронгауз Максим Анисимович	Институт Лингвистики РГГУ
Лукашевич Наталья Валентиновна	НИВЦ МГУ
Маккарти Диана	Lexical Computing Ltd., UK
Мельчук Игорь Александрович	Монреальский университет
Нивре Йоаким	Уппсальский университет
Ниренбург Сергей	Университет Нью-Мексико
Осипов Геннадий Семёнович	Институт программных систем РАН
Попов Эдуард Викторович	РосНИИ информационной техники

и САПР

Purdue University, USA

University of Leeds, UK

University of Southern California, USA

Компания Yandex

Компания АВВҮҮ

#### Организационный комитет

Селегей Владимир Павлович, Компания АВВҮҮ

председатель

Байтин Алексей Владимирович Компания Yandex

Беликов Владимир Иванович Институт русского языка им. В.В. Виноградова

PAH

Браславский Павел Исаакович Kontur Labs;

Уральский федеральный университет

Добров Борис Викторович НИВЦ МГУ

Иомдин Леонид Лейбович Институт проблем передачи информации РАН

Кобозева Ирина Михайловна Филологический факультет МГУ

Козеренко Елена Борисовна Институт проблем информатики РАН

Лауфер Наталия Исаевна ООО «проФан Продакшн»

Ляшевская Ольга Николаевна Universitet i Tromsø Сердюков Павел Викторович Компания Yandex

Соколова Елена Григорьевна РосНИИ искусственного интеллекта Толдова Светлана Юрьевна Филологический факультет МГУ

Шаров Сергей Александрович University of Leeds, UK

#### Секретариат

Белкина Александра Андреевна, секретарь оргкомитета Компания АВВҮҮ

Мытникова Татьяна Александровна, координатор Компания АВВҮҮ

#### Рецензенты

Августинова Таня Азарова Ирина Владимировна Апресян Валентина Юрьевна Байтин Алексей Владимирович Баранов Анатолий Николаевич Беликов Владимир Иванович Богданов Алексей Владимирович Богданова Наталья Викторовна Богуславский Игорь Михайлович Бонч-Осмоловская Анастасия Александровна Браславский Павел Исаакович Гельбух Александр Феликсович Горностай Татьяна Александровна Губин Максим Вадимович Даниэль Михаил Александрович Добров Борис Викторович Добровольский Дмитрий Олегович Добрынин Владимир Юрьевич Дружкин Константин Юрьевич Захаров Леонид Михайлович Иомдин Борис Леонидович Иомдин Леонид Лейбович Кибрик Андрей Александрович Кобозева Ирина Михайловна

Козеренко Елена Борисовна

Крейдлин Григорий Ефимович

Кронгауз Максим Анисимович

Кэрролл Джон Лахути Делир Гасемович Левонтина Ирина Борисовна Лобанов Борис Мефодьевич Лукашевич Наталья Валентиновна Ляшевская Ольга Николаевна Маккарти Диана Падучева Елена Викторовна Пазельская Анна Германовна Подлесская Вера Исааковна Савельев Василий Евгеньевич Селегей Владимир Павлович Семенова-Флюр Вера Эммануиловна Сердюков Павел Викторович Сокирко Алексей Викторович Соколова Елена Григорьевна Старостин Анатолий Сергеевич Тестелец Яков Георгиевич Тихомиров Илья Александрович Толдова Светлана Юрьевна Урысон Елена Владимировна Федорова Ольга Викторовна Филиппова Екатерина Александровна Хорошевский Владимир Федорович Циммерлинг Антон Владимирович Шаров Сергей Александрович Юдина Мария Владимировна Янко Татьяна Евгеньевна

#### Contents\*

#### Раздел I. Основная программа конференции

Akinina Y. S., Kuznetsov I. O., Toldova S. Y.
The impact of syntactic structure on verb-noun collocation extraction
Алпатов В. M.
Александр Евгеньевич Кибрик: от структурализма к новым идеям 17
Антонова А. Ю., Соловьев А. Н.
Использование метода условных случайных полей
для обработки текстов на русском языке
Апресян В. Ю.
Семантика эмоциональных каузативов:
статус каузативного компонента
Azimov A. E., Bolshakova E. I.
Correcting collocation errors in learners' writing based
on probability of syntactic links
Баранов А. Н.
Семантика угрозы в лингвистической экспертизе текста
Беликов В. И., Копылов Н. Ю., Пиперски А. Ч.,
Селегей В. П., Шаров С. А.
Корпус как язык: от масштабируемости
к дифференциальной полноте
Benigni V., Cotta Ramusino P.
Computational treatment of support verb constructions
in Italian and in Russian
Bocharov V. V., Alexeeva S. V., Granovsky D. V.,
Protopopova E. V., Stepanova M. E., Surikov A. V.
Crowdsourcing morphological annotation
Bogdanov A. V., Leontyev A. P.
Description of the Russian external possessor construction in a
natural language processing system 115

The reports of each section are ordered by the surname of the first author in compliance with the English alphabet.

Богданова-Бегларян Н. В. Кто ищет — всегда ли найдет? (о поисковой функции вербальных хезитативов в русской спонтанной речи)	125
Большакова Е. И., Большаков И. А. Компьютерный словарь русских паронимов, основанный на формальном критерии паронимии	137
Борисова Е. Г., Пирогова Ю. К. Моделирование нетривиальных условий понимания сообщения (на примере иронии)	148
Brykina M. M., Faynveyts A. V., Toldova S. Yu.  Dictionary-based ambiguity resolution in Russian named entities recognition. A case study	163
Chernyak E. L., Mirkin B. G.  Computational refining of a Russian-language taxonomy using Wikipedia	177
Даниэль М. А., Добрушина Н. Р. Русский язык в Дагестане: проблемы языковой интерференции	186
Дёгтева А. В., Азарова И. В. Структура эмоционально-экспрессивного компонента в тезаурусе русского языка RussNet	200
Dikonov V. G.  Development of lexical basis for the Universal Dictionary of UNL Concepts	212
Dobrovol'skij D. O.  German-Russian idioms online: on a new corpus-based dictionary	222
Федорова О. В., Деликишкина Е. А., Слабодкина Т. А., Ципенко А. А. Моделирование диалога в психолингвистике: взрослые и детские стратегии описания объектов действительности	230
Galitsky B., Ivovsky D., Kuznetsov S., Strok F.  Parse thicket representations of text paragraphs	239
Гилярова К. А. Статья такая статья. Об одном типе редупликации в современном русском языке	256
Grefenstette G. Linguistic analysis of social media	270
Гришина Е. А. Жестикуляционные профили русских приставок	271

Иомдин Л. Л.	
Читать не читал, но: об одной русской конструкции	
с повторяющимися словесными элементами	297
Иомдин Б. Л., Лопухина А. А., Панина М. Ф., Носырев Г. В., Вилл М. В., Зайдельман Л. Я., Матиссен-Рожкова В. И., Винокуров Ф. Г., Выборнова А. Н. <i>Маг вел мот</i> : изменения в языке на материале бытовой терминологии	311
Кашкин Е. В., Ляшевская О. Н. Семантические роли и сеть конструкций в системе FrameBank	325
Кибрик А. А. <b>Дискурсивная таксономия</b>	344
Киселева К. Л., Вознесенская М. М., Козеренко А. Д. Больше единицы: русские идиомы с компонентом один/един	345
Коротаев Н. А. Полипредикативные конструкции с то что в непубличной устной речи	358
Котов А. А. Компенсация коммуникативных стимулов в эмоциональном диалоге	368
Kotov A. A.  Compensation of communication stimuli in the emotional dialogue	368
Крейдлин Г. Е., Переверзева С. И. Тело и его части в разных языках и культурах (итоги научного проекта)	378
Кустова Г. И. Семантические механизмы формирования адвербиальных выражений на базе отглагольных существительных	392
Кюсева М. В., Резникова Т. И., Рыжова Д. А. Типологическая база данных адъективной лексики	407
Летучий А. Б.  Свойства нулевой связки в русском языке в сопоставлении со свойствами выраженного глагола	419
Левонтина И. Б. О причинном значении союза а то	434
Litvinenko A. O.  Reported speech in spoken discourse: intonation as a means of integration	446
Loginova-Clouet E. A., Daille B.  Multilingual compound splitting combining language dependent and independent features	455

Ляшевская О. Н., Митрофанова О. А., Паничева П. В.	
Визуализация данных для каталога русских лексических	
конструкций (на материале НКРЯ)	464
Ляшевская О. Н.	
частотный лексико-грамматический словарь: проспект проекта	478
The Total Market Control of the Cont	170
Микаэлян И. Л., Зализняк Анна А.	
Вместе или раздельно? Заметки о семантической категории	
парности в русском языке	490
Михеев М. Ю.	
Да черт ли в деталях? Мера для оценки совпадения элементов	
идиостиля в текстах одного — или двух разных? Авторов	
(Агеев — Сирин/Набоков — Леви)	504
(Teeb onpini, Indokob viebn)	501
Nedoluzhko A., Mírovský J., Novák M.	
A coreferentially annotated corpus and anaphora resolution for Czech	519
Nekhay I. V.	
The prospects of application of semantic markup to the named entity	
recognition problem	528
recognition problem	320
Падучева Е. В.	
Эгоцентрические единицы языка и режимы интерпретации	538
Панина М. Ф., Байтин А. В., Галинская И. Е.	
Автоматическое исправление опечаток в поисковых запросах	
без учета контекста	556
ocs y teru kontekeru	550
Paperno D. A., Roytberg A. M., Khachko D. V., Roytberg M. A.	
Breeds of cooccurrence: an attempt at classification	568
П А. Р.	
Пазельская А. Г. Инкорпорация в глагольных формах в русском языке	F70
инкорпорация в глагольных формах в русском языке	3/9
Пестова А. Р.	
Семантические факторы изменения управления	
существительных в современном русском языке	592
Пиперски А. Ч., Сомин А. А.	
Литуративы в русском интернете: семантика, синтаксис	
и технические особенности бытования	605
Подлесская В. И.	
Нечеткая номинация в русской разговорной речи:	
опыт корпусного исследования	619

Поляков А. Е., Савчук С. О., Сичинава Д. В.	
Грамматический словарь для автоматического анализа	
текстов XVIII–XIX века: первые результаты	632
Protopopova E. V., Bocharov V. V.	
Unsupervised learning of part-of-speech disambiguation rules	655
Рахилина Е. В.	
Кондуктор, нажми на тормоза	665
Савинич Л. В.	
Использование контраста и эмфазы для передачи имплицитных смыслов	674
Семенова С. Ю.	
О полисемии «параметр — большое значение параметра»	688
Шилихина К. М.	
Дискурсивное маркирование нетривиального лексического выбора	698
Skopinava A. M., Hetsevich Yu. S., Lobanov B. M.	
Processing of quantitative expressions with units of measurement in	
scientific texts as applied to Belarusian and Russian text-to-speech synthesis	708
Slioussar N. A., Cherepovskaia N. V.	
Processing of case morphology: evidence from Russian	726
Соколова Е. Г., Кононенко И. С.	
Какие «ситуации» обозначаются русскими глаголами	
«отличить — отличать»	736
Solovyev V. D., Polyakov V. N.	
Database "Languages of the World" and it's application. State of the art	748
Татевосов С. Г.	
Грамматика глагола и диалектное варьирование	759
Урысон Е. В.	
Неотрицаемые предикаты: наречие впору	772
Yanko T. E.	
Sentence incompleteness vs. Discourse incompleteness: pitch accents	
and accent placement	783
Zhila A., Gelbukh A.	
Comparison of open information extraction for English and Spanish	794
Zimmerling A. V.	
Transitive impersonals in Slavic and Germanic:	
zero subjects and thematic relations	803

## Раздел I.

Основная программа конференции

# ЧАСТОТНЫЙ ЛЕКСИКО-ГРАММАТИЧЕСКИЙ СЛОВАРЬ: ПРОСПЕКТ ПРОЕКТА<sup>1</sup>

Ляшевская О. H. (olesar@gmail.com)

НИУ Высшая школа экономики, Москва, Россия

Обсуждается задача создания электронного частотного словаря, в котором будет отражено распределение грамматических форм в парадигме словоизменения русских имен существительных, прилагательных и глаголов, т.е. грамматический профиль индивидуальных лексем и лексических групп. В практике составления частотных словарей и квантитативных исследований стандартным объектом изучения является общая иерархия грамматических категорий, например, частотность частеречных классов или среднее соотношение частот именительного и творительного падежей. В данном проекте фокус переносится на распределение грамматических форм у конкретных лексем, выявление единиц с нестандартным перевесом тех или иных форм в парадигме. Словарь предназначен для исследований русской грамматики, грамматической семантики, а также изучения вариативности форм.

Ресурс строится на материалах Национального корпуса русского языка. В статье затрагиваются общие вопросы использования корпусов для создания частотных ресурсов подобного рода и технологии обработки данных. Предлагаются решения, связанные с отбором данных, уровнем дробности грамматических кластеров, параметрами мониторинга изменения грамматического профиля в зависимости от времени создания текста и жанрово-функционального регистра.

**Ключевые слова:** частотный словарь, грамматический профиль лексемы, словоизменение, грамматическая семантика, вариативность, русский язык, НКРЯ

# LEXICO-GRAMMATICAL FREQUENCY DICTIONARY: A PRELIMINARY DESIGN

Lyashevskaya O. N. (olesar@gmail.com)

NRU Higher School of Economics, Moscow, Russia

A new electronic frequency dictionary shows the distribution of grammatical forms in the inflectional paradigm of Russian nouns, adjectives and verbs,

В работе использованы результаты, полученные в рамках проекта № 11-01-0171, выполненного в рамках Программы «Научный фонд НИУ ВШЭ» в 2012–2013 гг.

i.e. the grammatical profile of individual lexemes and lexical groups. While the frequency hierarchy of grammatical categories (e.g. the frequency of part of speech classes or the average ratio of Nominative to Instrumental case forms) has long been the standard topic of research, the present project shifts the focus to the distribution of grammatical forms in particular lexical units. Of particular concern are words with certain biases in grammatical profile, e. g. verbs used mostly in Imperative, in past neutral or nouns used often in plural. The dictionary will be a source for many of the future research in the area of Russian grammar, paradigm structure, grammatical semantics, as well as variation of grammatical forms.

The resource is based on the data of the Russian National Corpus. The article addresses some general issues such as corpora use in compiling frequency resources and technology of corpus data processing. We suggest certain solutions related to the selection of data and the level of granularity of grammatical profile. Text creation time and language registers are discussed as parameters which may shape the grammatical profile fluctuations.

**Key words:** frequency dictionary, grammatical profile, inflection, semantics of grammar, form variation, Russian, Russian National Corpus

#### 1. Введение

Частотный лексико-грамматический словарь продолжает серию частотных словарей, создаваемых на данных Национального корпуса русского языка, и является прямым продолжением частотного словаря (Ляшевская, Шаров 2009). В общем частотном словаре основная доля информации была представлена на уровне лексем. Из грамматической информации давались сведения о доле слов разных частей речи и о наиболее частотных словоформах русского языка. Вместе с тем, если смотреть с точки зрения конкретной лексемы, информации о частоте всех ее словоформ словарь не давал. Эту лакуну заполняет новый экспериментальный лексико-грамматический словарь. Он представляет грамматический профиль (т.е. распределение грамматических форм в парадигме словоизменения) 5000 наиболее частотных русских имен существительных, прилагательных и глаголов.

Далее в статье речь пойдет о задачах словаря, его структуре, а также о некоторых проблемных точках, связанных с обработкой И интерпретацией частотных данных.

#### 2. Предназначение словаря

Квантитативныеисследованиянелексических единицязыка—грамматических классов (например, иерархий падежного маркирования), грамматических форм внутри парадигмы конкретного слова, вариативности грамматических

и лексико-грамматических единиц, вариативности падежного и предложнопадежного оформления зависимых — были признаны необходимой составляющей лингвистического анализа еще в мировой лингвистике 50–70-х годов ХХ в. В русистике были получены замечательные результаты в классических работах Штейнфельдт 1963, Greenberg 1974, Граудина и др. 1976, Апресян 1967 и мн. др.). Однако именно появление представительных и сбалансированных лингвистических корпусов объемом от ста миллионов словоупотреблений и выше поставило эти исследования на принципиально новый уровень, как в плане используемых математических статистических моделей и компьютерных технологий, так и в плане осмысления частотных результатов и их устойчивости.

В теоретической лингвистике частотные исследования приобрели особую актуальность в связи с постулированием usage-based model — модели языка, предполагающей, что частота употребления языковых единиц оказывает непосредственное влияние на их конструктивные свойства, статус в системе, вариативность и изменение в истории языка (Kemmer & Barlow 2000). Еще одна гипотеза — о семантической мотивированности грамматических явлений верифицируется в ходе исследований, изучающих сдвиги частот форм в разных лексико-семантических классах (см. об этом Janda. Lvashevskava 2011): например, предполагается, что преобладание форм императива несовершенного вида связано с семантическими и функциональными особенностями лексических единиц. В когнитивных исследованиях изучается также гипотеза о том, что возможности языковой памяти таковы, что в частотных фрагментах человек оперирует единицами, большими чем слово (pre-fabricated units). Поднимается и вопрос, оперирует ли человек лексемами, т.е. единицами абстрактного уровня, или же это порождение грамматической схоластики, и человек опирается в своем языковом опыте исключительно на словоформы (Newman 2008). Наконец, изучение грамматических частотных профилей в разных языках помогло бы извлечь новые факты для лингвистической типологии и истории развития языков.

В грамматике русского языка, и теоретической, и практической, традиционно большую роль играет вопрос о дефектных парадигмах, а также о вариативных формах словоизменения. Несмотря на получившую общее признание точку зрения о градуальности таких явлений, как, например, singularia и pluralia tantum, выявление ассоциированных с ними лексических единиц и описание их функционирования все еще нуждается в эмпирических данных. То же можно сказать и о проблематике появления, со-существования и исчезновения вариативных форм типа род. мн. помидор/помидоров, прош. ед. стыл/стынул, статусе «вторых» падежей и т.п.

В преподавании родного и иностранных языков знание о частотных фактах грамматики позволяет выстроить правильную последовательность изучения грамматических тем (например, порядок изучения падежей), соотнести грамматические категории с теми лексемами, при которых они чаще всего встречаются, изучать лексику в контексте (знать самые частотные сочетания), выбирать для образца тексты, наиболее подходящие по жанрово-стилевому признаку к изучаемой грамматической теме и т. п.

И, конечно, неоценимую роль играют частотные данные в разработках систем автоматической обработки текста. Особенно это стало очевидно в эпоху стремительного развития алгоритмов машинного обучения, построенных на вероятностях. Грамматические и сочетаемостные преференции слов учитываются в синтаксических парсерах, системах разрешения неоднозначности, средствах исправления орфографии, распознавания текста, в онтологических расширениях поисковых систем и др.

Несмотря на то, что задача построения частотной русской грамматики и фронтального изучения грамматической вариативности осознана и ставится в литературе (Мустайоки 1973, Baerman et al. 2010), в настоящее время не существует ни одного сколько-нибудь полного лексикографического ресурса, приближающего нас к этой цели. Ресурс на материале НКРЯ дает уникальную возможность ответить на многие исследовательские вопросы, исходя из современных возможностей корпусной лингвистики.

# 3. Общая температура по больнице, или почему не всегда помогает статистика падежей

Когда говорят о частотной грамматике языка, в первую очередь, имеют в виду соотношения частот частеречных классов, падежей и других грамматических категорий. Особенно популярна тема частотного распределения падежей — в работе Копотев 2008 цитируются три исследования, появившихся только в 1959-1961 гг., что касается настоящего времени, то, как показывает веб-поиск, аналогичные работы, построенные на разных текстовых выборках, плодятся с невиданной скоростью. Работа самого М. Копотева привлекает внимание к устойчивости частотных данных на больших корпусах (см. табл. 1). Его вывод — в том, что современные корпуса довольно хорошо согласуются друг с другом в оценке средней вероятности появления падежей, а различия кроются в жанровой принадлежности текстов.

и в п							
падежей по данным (Копотев 2008)							
Табл. 1. Частотное распределение шести							

	И	P	Д	В	Т	П
□НКРЯ	27,06	29,23	5,98	18,66	8,44	10,63
■ XAHKO	24,30	32,62	5,50	17,73	8,08	11,78
□ J. 1953	38,80	16,80	4,70	26,30	6,50	6,90
<b>■</b> Št. 1963	33,60	24,60	5,10	19,50	7,80	9,40

Однако, легко видеть, что принцип «выбирай родительный, если забыл — не ошибешься» может сыграть злую шутку со студентом РКИ, в случае, если ему нужно употребить слово menom. Как показывает табл.  $2^2$ , распределение

 $<sup>^{2}~</sup>$  Здесь и далее в таблицах приведены данные по корпусу со снятой лексико-грамматической омонимией НКРЯ.

частот падежей у некоторых существительных может разительно отличаться от средней картины.

И	Р	Л	В	Т	П	Всего (
шепот,	поза, тр	опинка (і	падежнь	іе формі	ы ед. чис	сла)
					T	

Табл. 2. Частотный грамматический профиль лексем

	И	P	Д	В	T	П	Bcero (F.abs)
шепот	10,9%	3,7%	0,9%	8,3%	75,6%	0,6%	349
поза	15,9%	6,3%	0,8%	19,0%	4,0%	54,0%	126
тропинка	27,6%	2,0%	52,0%	5,1%	5,1%	8,2%	98

Дж. Гринбергу принадлежит наблюдение, что разные семантические группы должны иметь разную дистрибуцию падежей (как в предложных, так и в беспредложных употреблениях), иными словами, средние значения падежных показателей в группе имен абстрактных качеств (или имен частей тела, или названий мер) должны отличаться от средних значений по всему массиву лексики (Greenberg 1974/1991). Выбор русского языка как объекта исследования Гринберга был не случаен — именно в тот момент русский язык, один из немногих, располагал частотным списком форм падежей и предложно-падежных сочетаний имен существительных, входившим в состав замечательного частотного словаря Э. Штейнфельдт (Šteinfeldt 1963). Гринберг искал «волшебное» соотношение, которое позволяло бы отнести слово к тому или иному семантическому классу — и, естественно, не нашел его. Позднее его наблюдение было реинтерпретировано как семантически мотивированный сдвиг частот грамматических форм. Например, большую долю форм творительного падежа шепотом легко объяснить пересечением в семантике грамматической формы (творительный способа) и семантике лексемы (шепот как способ произнесения); форм предложного падежа (в) позе — связью между стативной семантикой существительного и семантикой локативной группы  $\beta + S.loc$ , наиболее типичном контекстном варианте употребления этого слова. Аналогичным образом, преобладание форм датива у существительного тропинка объясняется тем, что слова со значением траектории — идеальный лексический наполнитель предложной группы no + S.dat.

В работе (Janda, Lyashevskaya 2011) мы ввели понятие грамматического профиля лексемы — как инструмента для изучения семантических и функциональных причин девиаций грамматических форм. Исследование поведения форм вида, времени и наклонения, частности, показало предсказуемые частотные эффекты в разных клетках парадигмы: в императиве несовершенного вида — для глаголов привлечения внимания, вежливой просьбы, лексики, относящейся к культурному фрейму встречи гостей и т.п., ср. раздевайтесь, садитесь, присоединяйтесь, закусывайте, закуривайте, ступайте, прощайте; в инфинитиве совершенного вида — для глаголов, в которых заложена презумпция труднодостижимого результата (вследствие этого они часто употребляются в контексте глаголов попытки, модальных предикативов, в целевых придаточных и т.п., ср. попытался/тяжело было/чтобы восполнить) и т.п.

В исследовании (Kuznetsova 2013) вводятся классы типичных «женских» и типичных «мужских» глаголов — соотношение форм мужского и женского рода у глаголов типа вышивать и глаголов типа надвинуть будет разным.

На материале BNC С. Райс и Дж. Ньюман (Rice, Newman 2005, Newman 2008) сделали наблюдение, что разброс грамматического распределения может присутствовать и внутри лексических групп. Они показали, что даже близкие по смыслу слова, английские think, know и mean, могут иметь значительную диспропорцию форм времени, лица и числа, и назвали это явление "inflectional islands". Объяснение этого явления кроется в индивидуальных семантических особенностях каждого глагола, в способности присоединять разные типы субъектов и т.п. В (Janda, Lyashevskaya 2011) указывается также большой вклад устойчивых конструкций в формирование тех или иных грамматических «флюсов» у индивидуальных лексем, ср. мне плевать, мне наплевать, на чужой каравай рот не разевай, хоть залейся, поминай, как звали.

Однако, наиболее удивительный факт русской лексической системы состоит в том, что почти не существует существительных, грамматический профиль которых соответствовал бы «среднему» профилю нарицательной лексики, глаголов со «средней» пропорцией форм времени-лица-числа и т. п. По сути, мы имеем дело со сложным наслаиванием семантических особенностей, сочетаемостных и конструктивных свойств, которые суммарно влияют на частотный выход.

#### 4. Обработка корпусных данных

Основная часть словаря строится на данных 1900–2010 гг., в диахронической части привлекаются данные, начиная с 1800 г. Данные для «малого» словаря были собраны по корпусу со снятой лексико-грамматической омонимией (5,4 млн словоупотреблений, стандартная коллекция), для «большого» словаря — по основному, газетному, поэтическому и устному корпусу. Сбор осуществлялся с учетом функциональных стилей и жанров текста, а также с учетом времени создания.

Прежде всего, была собрана статистика по словоформам с лексико-грамматическим разбором (лемма, часть речи, словоизменительные характеристики)<sup>3</sup>, разметкой лексико-семантического класса капитализации написания. Были также собраны 2- и 3-грамы, отражающие статистику предложно-падежных сочетаний существительных и местоимений.

Для «борьбы» с грамматической омонимией словоформ внутри парадигм и между парадигмами использовалась автоматически дизамбигуированная версия основного, газетного, поэтического и устного корпуса. Она была создана с применением двух программ — модуля на эвристиках и НММ-модуля, обученного на текстах снятого вручную корпуса. Небольшая часть данных дополнительно корректировалась вручную.

Особо отметим, что большую проблему для дизамбигуации представляют ингерентно-пересеченные парадигмы, например, парадигмы мужского

<sup>&</sup>lt;sup>3</sup> Использовались стандартные соглашения словаря (Ляшевская, Шаров 2009).

и женского рода имени рояль или парадигмы прилагательных вида запасной и запасный. Устаревший вариант женского рода существительного распознается словарем лишь в формах, не предусмотренных в парадигме мужского рода (роялью), и тем самым, в словаре отражается искусственно дефектная парадигма. Пересеченные парадигмы прилагательных, различающихся лишь в именительном падеже, также разводятся плохо, поскольку модели дизамбигуации не предусматривают столь тонкой настройки, да и вручную в письменном корпусе далеко не всегда удается однозначно определить лексему. Такие точечные места в словаре, где информация может быть недостоверна по причине несовершенной дизамбигуации, снабжаются специальной пометой.

#### 5. Виды частотной информации в словаре

Пользователь имеет возможность пользоваться двумя наборами данных. «Малый» словарь представляет наиболее аккуратные результаты в смысле разведения омонимов. Однако в корпусе со снятой вручную омонимией многие многие грамматические формы частотных лексем могут быть либо не представлены вообще, либо встречаются редко, и следовательно, не могут показать достоверное распределение форм. «Большой» словарь строится на корпусах НКРЯ, в десятки раз превосходящих «снятник», однако следует учитывать, что в некоторых зонах (например, в зоне противопоставления родительного и винительного падежа одушевленных существительных) информация в нем менее достоверна.

#### 5.1. Грамматические категории

Пользователь может выбрать данные как по всем грамматическим формам парадигмы, так и по более крупным кластерам форм. Например, могут быть приведены суммарные данные по формам полных пассивных причастий (без учета признаков падежа, числа и рода), по четырем формам прошедшего времени глагола, по всем формам единственного VS множественного числа существительного. Информация о падежных распределениях существительных и местоимений дополнена сведениями о распределении предложных конструкций, в которых задействован тот или иной падеж. Кроме того, можно получить сопоставительные данные для написаний с прописной VS строчной буквы.

#### 5.2. Омонимия и вариативность

Из всей парадигмы пользователю могут быть выданы сведения только об омонимичных формах (в т.ч. внутрипарадигматическая омонимия, ср.

солдат — им. ед. и род. мн., омонимия форм, принадлежащих разным парадигмам, ср. заплыв — формы имени существительного и глагола, см. Венцов, Касевич 2004). Предоставляются сведения о соотношении частот вариантов грамматических форм (например, сильней и сильнее, дверями и дверьми), так наз. «основных» и «вторых» падежей, различающихся на письме (ср. без толка и без толку), и других секундарных форм (ср. сильней и посильней).

#### 5.3. Распределение по годам и жанрам

Информация об изменении грамматических профилей во времени дается в 10-летних интервалах; в газетном корпусе учитываются интервалы в 1 год. Пользователь может увидеть распределения в художественной прозе, в поэзии, в периодике, в бытовой, учебно-научной и т.п. сферах нехудожественной литературы, в электронной коммуникации, а также в устной непубличной речи.

#### 5.4. Единицы измерения

Пользователь может выбрать один или несколько вариантов представления частотной информации:

- количество текстов корпуса, в которых встретились формы;
- абсолютная частота вхождений и размер корпуса;
- частота в ірт;
- иерархия форм у рассматриваемой единицы/класса вида Loc > Gen > Nom > Acc > Dat > Ins;
- процентное распределение (см. табл. 3) и попарное соотношение форм;
- квинтильное распределение каждой из форм, например, положение формы предложного падежа единственного числа слова *велосипед* в первой, второй... пятой порции списка, в котором представлены формы предложного падежа единственного числа всех существительных (а самые редкие, д самые частые, см. табл. 4).

**Табл. 3.** Профиль падежных форм лексемы *влияние*: абсолютное и относительное распределение

	И	P	Д	В	T	П	Bcero (F.abs)
sg	98	128	29	170	137	14	576
pl	4	9	3	7	2	2	27
				_	_	_	_
	И	P	Д	В	T	П	Bcero (%)
sg	И 17,0%	P 22,2%	Д 5,0%	B 29,5%	T 23,8%	П 2,4%	Bcero (%) 100,0%

Табл. 4. Квинтильное распределение падежных форм
ед. числа в группе имен транспортных средств

Лемма	И	P	Д	В	Т	П	Всего (F.abs)
метро	a	Д	Г	a	a	Д	185
корабль	Д	В	б	б	a	В	231
грузовик	Д	Г	В	б	б	В	134
пароход	Д	Д	a	б	В	Г	121
автомобиль	Г	Г	В	б	б	Г	441
поезд	Д	В	В	б	б	Г	618
самолет	Г	В	Г	В	В	Г	385
трамвай	Г	б	В	Г	В	Г	198
лодка	Г	В	б	Г	б	Г	280
вагон	a	Г	Г	В	a	Д	473
велосипед	б	В	a	Г	б	Д	206
автобус	Г	б	В	В	б	Д	281
такси	В	a	б	Д	a	Д	174

Оговорим, что пользователь может выбрать разные методики расчета соотношений частот в парадигме, известных из литературы. За основу сравнения (100%) может быть принята вся парадигма (т.е. сумма всех частот грамматических форм), некоторая базовая часть (например, парадигма глагола за вычетом форм причастий и деепричастий), приоритетная форма (например, сумма форм прошедшего времени), а также доля употреблений двух форм относительно друг друга (например, отношение частоты форм женского рода к частоте форм мужского рода).

#### 5.5. Сравнение лексем. Классы

Информация в словаре разнесена на несколько уровней. Первый уровень — индивидуальные грамматические профили лексем. На втором уровне даются сведения для крупных лексико-семантических классов (в классификации НКРЯ), например, для глаголов движения, имен инструментов и т.п. Третий уровень — распределение грамматических частот на уровне частеречного класса (словарь также дает справочную информацию о встречаемости самих частеречных классов, а также именных и глагольных грамматических категорий).

Таким образом, информация об индивидуальных лексемах может быть сопоставлена с данными по их лексико-семантическому классу и, шире, со средним грамматическим профилем части речи. Возможно сопоставление грамматических профилей нескольких лексем между собой.

#### 6. Заключение

Словарь адресован, в первую очередь, исследователям русского словоизменения, грамматической семантики, тем, кто изучает вариативность грамматической нормы. Вместе с тем, нужно заметить, что «лексикоцентричный» подход, несмотря на ресурсоемкость и неплотность данных, может оправдывать себя и в автоматической обработке текста. В частности, в экспериментах (Данилова и др. 2013) показано, что учет лексического фактора позволяет повысить качество автоматической дизамбигуации лексико-грамматической омонимии на 3%.

Электронная форма словаря позволяет постоянно совершенствовать его. Во-первых, планируется развивать функционал с учетом пожеланий пользователей, в частности, дополнить словарь модулем графического представления результатов, подключить внешние словари (словарь вариантов, словообразовательный и т.п.) и др. Во-вторых, будет совершенствоваться качество даных за счет улучшения дизамбигуации корпусных данных и работы с сообщениями пользователей об ошибках. В-третьих, увеличение объема словаря: включение новых лексических данных, добавление информации об авторе и т.п., — требует дополнительных исследований, поскольку работа с малыми частотами (sparse data) требует особой осторожности и особых техник.

Главный вопрос — в том, как интерпретировать полученные данные, каким образом переносить сведения о статистических вероятностях на другие текстовые корпуса и как научиться делать аккуратные выводы о функционировании грамматических форм в целом. Предлагаемый словарь — лишь первый опыт составления большого лексико-грамматического ресурса, и, соответственно, станет благодатным материалом для исследования достоверности корпусных данных. Безусловно, мы должны лучше понимать структуру выборок, как она связана с устойчивостью статистических данных, научиться балансировать выборки для разных временных срезов, провести множество экспериментов с полученным лексическим материалом для того, чтобы достоверность интерпретации корпусных данных перестала вызывать вопросы.

#### Литература

- 1. Baerman M., Brown D., Corbett G. G., Krasovitsky A., Williams P. (2010), Predicate agreement in Russian: A corpus-base approach, Wiener Slawistischer Almanach, Sonderband 74, pp.109–121.
- 2. *Greenberg J. H.* (1974/1990), The relation of frequency to semantic feature in a case language (Russian), in Denning K., Kemmer S. (eds), On Language, Selected Writings of Joseph H. Greenberg, Stanford, pp. 207–226.
- 3. *Ilola E., Mustajoki A.* (1989), Report on Russian Morphology as it Appears in Zaliznyak's Grammatical Dictionary, (Slavica Helsingiensia 7), Helsinki.
- 4. *Janda L. A., Lyashevskaya O.* (2011), Grammatical profiles and the interaction of the lexicon with aspect, tense and mood in Russian, Cognitive Linguistics, 22 (4), pp. 719–763.
- 5. Kemmer S., Barlow M. (2000), A Usage-Based Conception of Language, Essen, 2000.
- 6. *Kuznetsova J.* (2013), Linguistic Profiles: Correlations between Form and Meaning. Ph.D. diss., University of Tromsø.
- Newman J. (2008), Aiming low in linguistics: Low-level generalizations in corpus based research. Proceedings of the 11th International Symposium on Chinese Languages and Linguistics, National Chiao Tung University, Hsinchu, Taiwan, May 24, 2008.
- 8. *Rice S., Newman J.* (2005), Inflectional islands, ICLC-9, Yonsei University, Seoul, Korea.
- 9. Šteinfeldt E. (1963), Russian Word Count, Moscow.
- 10. *Апресян Ю. Д.* (1967), Экспериментальное исследование семантики русского глагола, М.
- 11. *Венцов А. В., Касевич В. Б.* (ред.) (2004), Словарь омографов русского языка, СПб.: Филологич. ф-т СПбГУ.
- 12. Граудина Л. К., Ицкович В. А., Катлинская Л. П. (1976), Грамматическая правильность русской речи. Стилистический словарь вариантов. М.
- 13. Данилова В., Волков О., Ладыгина А., Привознов Д., Сербинова И., Сим Г. (2013). Снятие омонимии методом НММ (рукопись).
- 14. *Копотев М.* (2008), К построению частотной грамматики русского языка: падежная система по корпусным данным // Мустайоки А., Копотев М. В., Бирюлин Л. А., Протасова Е. Ю. (ред.), Инструментарий русистики: корпусные подходы, Хельсинки.
- 15. *Ляшевская О. Н., Шаров С. А.* (2009), Частотный словарь современного русского языка (на материале Национального корпуса русского языка), М.: Азбуковник.
- 16. *Мустайоки А.* (1973), Опыт составления частотной грамматики русских существительных, Хельсинки, (рукопись).

#### References

- 1. *Apresjan Ju. D.* (1967), Experimental research on the semantics of the Russian verb [Eksperimental'noe issledovanie semantiki russkogo glagola], Moscow.
- 2. Baerman M., Brown D., Corbett G. G., Krasovitsky A., Williams P. (2010), Predicate agreement in Russian: A corpus-base approach, Wiener Slawistischer Almanach, Sonderband 74, pp. 109–121.
- 3. *Danilova V., Volkov O., Ladygina A., Privoznov D., Serbinova I., Sim G.* (2013). Disambiguation with HMM [Snjatie omonimii metodom HMM] (manuscript).
- 4. *Graudina L. K., Ickovich V. A., Katlinskaja L. P.* (1976), Correct Russian speech: Stylistical dictionary of grammatical choices [Grammaticheskaja pravil'nost' russkoy rechi. Stilisticheskiy slovar' variantov]. Moscow.
- 5. *Greenberg J. H.* (1974/1990), The relation of frequency to semantic feature in a case language (Russian), in Denning K., Kemmer S. (eds), On Language, Selected Writings of Joseph H. Greenberg, Stanford, pp. 207–226.
- 6. *Ilola E., Mustajoki A.* (1989), Report on Russian Morphology as it Appears in Zaliznyak's Grammatical Dictionary, (Slavica Helsingiensia 7), Helsinki.
- 7. *Janda L. A., Lyashevskaya O.* (2011), Grammatical profiles and the interaction of the lexicon with aspect, tense and mood in Russian, Cognitive Linguistics, 22 (4), pp. 719–763.
- 8. Kemmer S., Barlow M. (2000), A Usage-Based Conception of Language, Essen, 2000.
- 9. *Kopotev M.* (2008), Towards the frequency grammar of Russian: corpus evidence on the grammatical case system [K postroeniju chastotnoy grammatiki russkogo jazyka: padezhnaja sistema po korpusnym dannym] // Mustayoki A., Kopotev M. V., Birjulin L. A., Protasova E. Ju. (eds.), Instruments of Russian linguistics: corpus approaches [Instrumentariy rusistiki: korpusnye podkhody], Helsinki.
- 10. *Kuznetsova J.* (2013), Linguistic Profiles: Correlations between Form and Meaning. Ph.D. diss., University of Tromsø.
- 11. *Lyashevskaya O. N., Sharoff S. A.* (2009), Frequency dictionary of modern Russian based on the Russian National Corpus [Chastotnyy slovar' sovremennogo russkogo jazyka (na materiale Nacional'nogo korpusa russkogo jazyka)], Azbukovnik, Moscow.
- 12. *Mustajoki A.* (1973), On compiling the frequency dictionary of Russian nouns [Opyt sostavlenija chastotnoy grammatiki russkikh suschestvitel'nykh], Hel'sinki, (manuscript).
- 13. *Newman J.* (2008), Aiming low in linguistics: Low-level generalizations in corpus based research. Proceedings of the 11th International Symposium on Chinese Languages and Linguistics, National Chiao Tung University, Hsinchu, Taiwan, May 24, 2008.
- 14. *Rice S., Newman J.* (2005), Inflectional islands, ICLC-9, Yonsei University, Seoul, Korea
- 15. Šteinfeldt E. (1963), Russian Word Count, Moscow.
- 16. *Ventsov A. V., Kasevich V. B.* (eds.) (2004), Dictionary of Russian homographs [Slovar' omografov russkogo jazyka], St.-Petersburg.

cesses corpus samples and learns up the constructions is described. We provide analysis for the structure and content of extracted constructions (e. g. r:ord der:num t:ord r:qual|pervyj 'first' + LJUBOV' 'love'; LJUBOV' 'love' + PR|s 'from' + ANUM m sg gen|pervyj 'first' + S f inan sg gen|vzgljad 'sight' = love at first sight). As regards their structure, constructions may be considered as n-grams (n is 2 to 5). The representation of constructions is bipartite as they may combine either morphological and lemma tags or lexical-semantic and lemma tags. We discuss the use of visualization module PATTERN.GRAPH that represents the inner structure of extracted constructions.

#### LEXICO-GRAMMATICAL FREQUENCY DICTIONARY: A PRELIMINARY DESIGN

Lyashevskaya O. N. (olesar@gmail.com), NRU Higher School of Economics, Moscow, Russia

A new electronic frequency dictionary shows the distribution of grammatical forms in the inflectional paradigm of Russian nouns, adjectives and verbs, i.e. the grammatical profile of individual lexemes and lexical groups. While the frequency hierarchy of grammatical categories (e.g. the frequency of part of speech classes or the average ratio of Nominative to Instrumental case forms) has long been the standard topic of research, the present project shifts the focus to the distribution of grammatical forms in particular lexical units. Of particular concern are words with certain biases in grammatical profile, e.g. verbs used mostly in Imperative, in past neutral or nouns used often in plural. The dictionary will be a source for many of the future research in the area of Russian grammar, paradigm structure, grammatical semantics, as well as variation of grammatical forms.

The resource is based on the data of the Russian National Corpus. The article addresses some general issues such as corpora use in compiling frequency resources and technology of corpus data processing. We suggest certain solutions related to the selection of data and the level of granularity of grammatical profile. Text creation time and language registers are discussed as parameters which may shape the grammatical profile fluctuations.

### TOGETHER OR SEPARATELY? ON THE SEMANATIC CATEGORY OF TWONESS IN RUSSIAN

**Mikaelian I.** (irina-mikaelian@yandex.ru), The Pennsylvania State University, State College, PA, USA, **Zalizaniak Anna A.** (anna.zalizaniak@gmail.com), Institute of Linguistics, Russian Academy of Sciences: Institute of Informatics. Russian Academy of Sciences. Moscowm Russia

This paper attempts to refine our understanding of the grammatical and semantic features of the Russian collective numerals using data of corpora. The focus of our attention is the word *dvoe* considered in comparison with other quantity words comprising the meaning 'two' in their semantics, i. e. the numerals *dva* 'two' and *oba* 'both', as well as the noun *para* 'pair, couple'. The importance for the Russian language of the semantic category of "twoness" has been shown, and a new term *gemina tantum* has been introduced to designate the class of nouns that tend to be used in plural form and normally refer to two objects forming a pair or a couple, cf. *shoes, boots, eyes, parents, spouses.* Semantic analysis of the words *dvoe* and *oba* in the context of human nouns has shown that these words practically never interchange because, despite similar assertions, they carry different presuppositions and implications.

# IS THE DEVIL IN THE DETAILS?... A MEASURE TO ASSESS THE MATCHING OF IDIOSTYLE ELEMENTS IN THE TEXTS OF ONE — OR IS IT TWO? — AUTHORS (AGEEV-SIRIN/NABOKOV-LEVI)

Mikheev M. Yu. (m-miheev@rambler.ru), NIVC MGU, Moscow, Russia

We analyze N. Struve's hypothesis that the author of the text of *The romance with cocaine* (published in 1936 under the pseudonym M. Ageev) was Vladimir Nabokov. We compare the idiostyle features of this text and all of Nabokov's texts, as well as what is available in the Russian National Corpus, published before the Ageev and Nabokov works and after them. The general conclusion is that Nabokov seems not to be involved in this text. This problem was stated by Nikita Struve, who rejected biographical arguments and required that "philological", literary or poetic arguments should be given. We consider all of these arguments.

### Авторский указатель

A H D 1 000	II
Азарова И. В т. 1: 200	Иомдин Б. Л т. 1: 311
Азимов А. Е т. 1: 61	Иомдин Л. Л т. 1: 297; т. 2: 132
Акинина Ю. С т. 1: 2	Кашкин Е. В т. 1: 325
Алексеева С. В т. 1: 109	Кибрик А. А т. 1: 344
Алпатов В. М т. 1: 17	Киселева К. Л т. 1: 345
Антонова А. Ю т. 1: 27	Клековкина М. В т. 2: 51
Апресян В. Ю т. 1: 44	Козеренко А. Д т. 1: 345
Байтин А. В т. 1: 556	Кононенко И. С т. 1: 736
Баранов А. Н т. 1: 72	Копылов Н. Ю т. 1: 83
Беликов В. И т. 1: 83	Корольков Е. А т. 2: 2
Белобородов А т. 2: 122	Коротаев Н. А т. 1: 358
Блинов П. Д т. 2: 51	Котельников Е. В т. 2: 51
Богданова-Бегларян Н. В т. 1: 125	Котов А. А т. 1: 368
Богданов А. В т. 1: 115	Крейдлин Г. Е т. 1: 378
Богуславский И. М т. 2: 132	Кузнецов И. О т. 1: 2
Большакова Е. И т. 1: 61, 137	Кузнецов С т. 1: 239
Большаков И. А т. 1: 137	Кузнецова Е. С т. 2: 71
Борисова Е. Г т. 1: 148	Кустова Г. И т. 1: 392
Бочаров В. В т. 1: 109, 655	Кюсева М. В т. 1: 407
Браславский П т. 2: 122	Левонтина И. Б т. 1: 434
Брыкина М. М т. 1: 163	Леонтьев А. Р т. 1: 115
Вилл М. В т. 1: 311	Летучий А. Б т. 1: 419
Винокуров Ф. Г т. 1: 311	Литвиненко А. О т. 1: 446
Вознесенская М. М т. 1: 345	Лобанов В. М т. 1: 708
Выборнова А. Н т. 1: 311	Логинова-Клуэ Е. А т. 1: 455
Галинская И. Е т. 1: 556; т. 2: 154	Лопухина А. А т. 1: 311
Галицкий Б т. 1: 239	Лукашевич Н. В т. 2: , 40
Гилярова К. А т. 1: 256	Людовик Т. В т. 2: 20
Грановский Д. В т. 1: 109	Ляшевская О. Н т. 1: 325, 464, 478
Гришина Е. А т. 1: 271	Мавлижутов Р. Р.
Гусев В. Ю т. 2: 154	Макеев И. В т. 2: 81
Даниэль М. А т. 1: 186	Марчук А. А т. 2: 81
Дёгтева А. В т. 1: 200	Матиссен-Рожкова В. И т. 1: 311
Деликишкина Е. А т. 1: 230	Мещерякова Е. М т. 2: 154
Диконов В. Г т. 1: 212; т. 2: 132	Микаэлян И. Л т. 1: 490
Добровольский Д. О т. 1: 222	Миркин Б. Г т. 1: 177
Добрушина Н. Р т. 1: 186	Митрофанова О. А т. 1: 464
Евдокимов Л. В т. 2: 145	Михеев М. Ю т. 1: 504
Зайдельман Л. Я т. 1: 311	Молчанов А. П т. 2: 145
Зализняк Анна А т. 1: 490	Нехай И. В т. 1: 528
Зуев К. А т. 2: 175	Носырев Г. В т. 1: 311
Иворский Д т. 1: 239	Остапук Н. А т. 2: 91
Инденбом Е. М т. 2: 175	Падучева Е. В т. 1: 538
	-

Пазельская А. Г т. 1: 579	Соловьев А. Н т. 1: 27
Панина М. Ф т. 1: 311, 556	Соловьев В. Д т. 1: 748
Паничева П. В т. 1: 464; т. 2: 101	Соломенник А. И т. 2: 31
· · · · · · · · · · · · · · · · · · ·	Сомин А. А т. 1: 605
Паперно Д. А т. 1: 568	
Переверзева С. И т. 1: 378	Степанова М. Е т. 1: 109
Пестов О. А т. 2: 51	Строк Ф т. 1: 239
Пестова А. Р т. 1: 592	Суриков А. В т. 1: 109
Пилипенко В. В т. 2: 20	Таланов А. О т. 2: 2
Пиперски А. Ч т. 1: 83, 605	Татевосов С. Г т. 1: 759
Пирогова Ю. К т. 1: 148	Тимошенко С. П т. 2: 132
Плешко В. В т. 2: 62	Толдова С. Ю т. 1: 2, 163
Подлесская В. И т. 1: 619	Уланов А. В т. 2: 81, 165
Поляков А. Е т. 1: 632	Урысон Е. В т. 1: 772
Поляков В. Н т. 1: 748	Федорова О. В т. 1: 230
Поляков П. Ю т. 2: 62	Фролов А. В т. 2: 62
Протопопова Е. В т. 1: 109, 655	Халилов М т. 2: 122
Рахилина Е. В т. 1: 665	Хачко Д. В т. 1: 568
Резникова Т. И т. 1: 407	Хетцевич Ю. С т. 1: 708
Ройтберг А. М т. 1: 568	Хомицевич О. Г т. 2: 11
Ройтберг М. А т. 1: 568	Циммерлинг А. В т. 1: 803
Рыжова Д. А т. 1: 407	Ципенко A. A т. 1: 230
Савинич Л. В т. 1: 674	Четверкин И. И т. 2:, 40
Савчук С. О т. 1: 632	Череповская Н. В т. 1: 726
Сапожников Г. А т. 2: 165	Черняк Е. Л т. 1: 177
Селегей В. П т. 1: 83	Чистиков П. Г т. 2: 2, 11, 31
Семенова С. Ю т. 1: 688	Чугреев А. А т. 2: 81
Сичинава Д. В т. 1: 632	Шаров С. А т. 1: 83; т. 2: 122
Скопинава А. М т. 1: 708	Шилихина К. М т. 1: 698
Слабодкина Т. А т. 1: 230	Шматова М. С т. 2: 154
Слюсарь Н. А т. 1: 726	Юдина М. В т. 2: 175
Соколова Е. Г т. 1: 736	Янко Т. Е т. 1: 783

#### **Author Index**

Akinina Y. S v. 1: 2	Granovsky D. V v. 1: 109
Alexeeva S. V v. 1: 109	Grefenstette G v. 1: 270
Alpatov V. M v. 1: 17	Grishina E. A v. 1: 271
Antonova A. Y v. 1: 27	Gusev V. Yu v. 2: 154
Apresjan V. Yu v. 1: 45	Hetsevich Yu. S v. 1: 708
Azarova I. V v. 1: 200	Indenbom E. M v. 2: 175
Azimov A. E v. 1: 61	Iomdin B. L v. 1: 312
Baranov A. N v. 1: 72	Iomdin L. L v. 1: 297; v. 2: 132
Baytin A. V v. 1: 556	Ivovsky D v. 1: 239
Belikov V v. 1: 84	Khachko D. V v. 1: 568
Beloborodov A v. 2: 122	Khalilov M v. 2: 122
Benigni V v. 1: 96	Khomitsevich O. G v. 2: 11
Blinov P. D v. 2: 51	Kiseleva K. L v. 1: 345
Bocharov V. V v. 1: 109, 655	Klekovkina M. V v. 2: 51
Bogdanova-Beglarian N. V v. 1: 125	Kononenko I. S v. 1: 736
Bogdanov A. V v. 1: 115	Kopylov N v. 1: 84
Boguslavsky I. M v. 2: 132	Korolkov E. A v. 2: 2
Bolshakova E. I v. 1: 61, 137	Korotaev N. A v. 1: 358
Bolshakov I. A v. 1: 137	Kotelnikov E. V v. 2: 51
Borisova E. G v. 1: 148	Kotov A. A v. 1: 368
Braslavski P v. 2: 122	Kozerenko A. D v. 1: 345
Brykina M. M v. 1: 163	Krejdlin G. E v. 1: 378
Cherepovskaia N. V v. 1: 726	Kustova G. I v. 1: 392
Chernyak E. L v. 1: 177	Kuznetsov I. O v. 1: 2
Chetviorkin I. I v. 2: 40, 71	Kuznetsov S v. 1: 239
Chistikov P. G v. 2: 2, 11, 31	Kuznetsova E. S v. 2: 71
Chugreev A. A v. 2: 81	Kyuseva M. V v. 1: 407
Cotta Ramusino P v. 1: 96	Leontyev A. P v. 1: 115
Daille B v. 1: 455	Letuchiy A. B v. 1: 420
Daniel M. A v. 1: 186	Levontina I. B v. 1: 434
Degteva A. V v. 1: 200	Litvinenko A. O v. 1: 446
Delikishkina E. A v. 1: 230	Lobanov B. M v. 1: 708
Dikonov V. G v. 1: 212; v. 2: 132	Loginova-Clouet E. A v. 1: 455
Dobrovol'skij D. O v. 1: 222	Lopukhina A. A v. 1: 312
Dobrushina N. R v. 1: 186	Loukachevitch N. V v. 2: 40, 71
Evdokimov L. V v. 2: 145	Lyashevskaya O. N v. 1: 465, 478
Faynveyts A. V v. 1: 163	Lyudovyk T. V v. 2: 20
Fedorova O. V v. 1: 230	Makeev I. V v. 2: 81
Frolov A. V v. 2: 62	Marchuk A. A v. 2: 81
Galinskaya I. E v. 1: 556; v. 2: 154	Màrquez L v. 2: 114
Galitsky B v. 1: 239	Matissen-Rozhkova V. I v. 1: 312
Gelbukh A v. 1: 794	Mavljutov R. R v. 2: 91
Gilyarova K. A v. 1: 256	Mescheryakova E. M v. 2: 154
•	•

Mikaelian I. L v. 1: 490	Savinitch L. V v. 1: 674
Mikheev M. Yu v. 1: 504	Selegey V v. 1: 84
Mirkin B. G v. 1: 177	Sharoff S v. 1: 84; v. 2: 122
Mírovský J v. 1: 519	Shilikhina K. M v. 1: 698
Mitrofanova O. A v. 1: 465	Sitchinava D. V v. 1: 633
Molchanov A. P v. 2: 145	Skopinava A. M v. 1: 708
Nedoluzhko A v. 1: 519	Slabodkina T. A v. 1: 230
Nekhay I. V v. 1: 528	Slioussar N. A v. 1: 726
Nosyrev G. V v. 1: 312	Shmatova M. S v. 2: 154
Novák M v. 1: 519	Sokolova E. G v. 1: 736
Ostapuk N. A v. 2: 91	Solomennik A. I v. 2: 31
Paducheva E. V v. 1: 538	Soloviev A. N v. 1: 27
Panicheva P. V v. 1: 465; v. 2: 101	Solovyev V. D v. 1: 748
Panina M. F v. 1: 311, 556	Somin A. A v. 1: 605
Paperno D. A v. 1: 568	Stepanova M. E v. 1: 109
Pereverzeva S. I v. 1: 378	Strok F v. 1: 239
Pestov O. A v. 2: 51	Surikov A. V v. 1: 109
Pestova A. R v. 1: 592	Talanov A. O v. 2: 2
Pleshko V. V v. 2: 62	Tatevosov S. G v. 1: 759
Piperski A v. 1: 84	Timoshenko S. P v. 2: 132
Piperski A. Ch v. 1: 605	Toldova S. Yu v. 1: 2, 163
Pirogova Yu. K v. 1: 148	Tsipenko A. A v. 1: 230
Podlesskaya V. I v. 1: 619	Ulanov A. V v. 2: 81, 165
Polyakov A. E v. 1: 633	Uryson E. V v. 1: 772
Polyakov P. Yu v. 2: 62	Vill M. V v. 1: 312
Polyakov V. N v. 1: 748	Vinokurov F. G v. 1: 312
Protopopova E. V v. 1: 109, 655	Voznesenskaja M. M v. 1: 345
Pylypenko V. V v. 2: 20	Vybornova A. N v. 1: 312
Rakhilina E. V v. 1: 665	Yanko T. E v. 1: 783
Reznikova T. I v. 1: 407	Yudina M. V v. 2: 175
Roytberg A. M v. 1: 568	Zajdel'man L. Ja v. 1: 312
Roytberg M. A v. 1: 568	Zalizaniak Anna A v. 1: 490
Ryzhova D. A v. 1: 407	Zhila A v. 1: 794
Sapozhnikov G. A v. 2: 165	Zuyev K. A v. 2: 175
Savchuk S. O v. 1: 633	

### Компьютерная лингвистика и интеллектуальные технологии

По материалам ежегодной Международной конференции «Диалог»

Выпуск 12 (19). 2013

Том 1. Основная программа конференции

Ответственный за выпуск **А. А. Белкина** Вёрстка **К. А. Климентовский** 

Подписано в печать 14.05.2013 Формат  $152 \times 235$  Бумага офсетная Тираж 250 экз. Заказ  $\mathbb{N}^{\circ}$  553

Издательский центр «Российский государственный гуманитарный университет» 125993, Москва, Миусская пл., д. 6 Тел.: +7 499 973 42 06

Отпечатано с готового оригинал-макета в типографии ООО «Издательско-полиграфический центр Маска» 117246, Москва, Научный пр-д, д. 20, стр. 9



Чтобы создавать современные поисковые технологии, нужны глубокие знания в области математики, лингвистики, анализа данных, программирования и других дисциплин. Благодаря давним традициям российской науки и образования, Яндексу удалось создать команду сильных специалистов, которые и сделали Яндекс одной из ведущих ІТ-компаний в России. Присылайте свое резюме, если хотите присоединиться к команде.

Все вакансии Яндекса: <a href="http://company.yandex.ru/job/vacancies/">http://company.yandex.ru/job/vacancies/</a>

#### Аналитик отдела веб-поиска

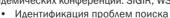
В обязанности входит анализ веб-страниц и запросов пользователей, а также их поискового поведения. Необходимо отличать предпочтения пользователей и качество поиска от статистического шума и артефактов. Стать аналитиком может человек с высшим техническим или математическим образованием, свободно читающий профессиональную литературу на английском. Пригодится также общее понимание поисковых технологий.



Пройти тестовое задание: http://clck.ru/AnX5

## Прикладной исследователь по направлению Data Mining / Information Retrieval

Работа для начинающих и опытных исследователей (data scientists, applied researchers). В обязанности входит разработка новых методов обработки информации, повышающих качество поиска, и описание их в статьях уровня ведущих академических конференций: SIGIR, WSDM, CIKM и т. п.



- Анализ существующих решений
- Проведение экспериментов
- Написание научных статей



Пройти тестовое задание: http://clck.ru/8deUG

#### Стажировка в Яндексе

Мы приглашаем студентов, аспирантов и выпускников вузов в московский и петербургский офисы Яндекса. Вы сможете своими глазами увидеть, как в Яндексе создают интернет-сервисы, поработать над «боевыми» задачами, поучиться у знатоков своего дела. Проявите себя, и, вполне возможно, мы пригласим вас на постоянную работу.

- Опыт работы не важен
- Работу в Яндексе можно совмещать с учёбой
- Стажёры получают зарплату и бесплатные обеды

Заполните анкету: <a href="http://company.yandex.ru/job/intern/">http://company.yandex.ru/job/intern/</a>