

СОЧЕТАЕМОСТЬ ЧЕРЕЗ ПРИЗМУ КОРПУСОВ

Захаров В. П. (vz1311@yandex.ru)

Санкт-Петербургский государственный
университет, Санкт-Петербург, Россия

Изучение сочетаемости слов является одной из главных задач лингвистики. Явное отражение это явление нашло в выделении устойчивых сочетаний, которые являются главным объектом фразеологии, и в создании словарей устойчивых словосочетаний. В статье рассматриваются устойчивые сочетания разного типа и показываются способы их количественной оценки. Описаны эксперименты, в ходе которых на материале корпусов русского языка и инструментов корпусной лингвистики, таких как Sketch Engine и Google books Ngram Viewer, было показано, как с помощью корпусных методов можно расширить состав словарных статей в словарях устойчивых выражений и как можно количественно оценить употребительность и устойчивость словосочетаний в синхронии и диахронии.

Ключевые слова: устойчивые словосочетания, фразеологизмы, коллокации, словари сочетаемости, корпусы текстов, меры ассоциации, диахронические исследования

SET PHRASES: A VIEW THROUGH CORPORA

Zakharov V. P. (vz1311@yandex.ru)

Saint-Petersburg State University, Saint-Petersburg, Russia

The study of word collocability is one of the main tasks of linguistics. Syntagmatic relations bind together language units being in direct contact with each other. The combinatory ability of language units, collocability, is one of the linguistic syntagmatic laws. This phenomenon is the main object of the phraseology and lexicography. The article deals with set phrases of different types from the point of view of their numerical evaluation. Corpus linguistics understand set phrases as statistically determined unities. This approach is the basic point of different automatic ways to extract idioms as collocations. The paper describes experiments which show how text corpora and corpus methods and tools such as association measures, word sketches, concordances can be used to expand the entries in existing dictionaries and how set phrases could be evaluated quantitatively. There are a small numbers of works on set phrases productivity during time periods because of small size of historical corpora. In this research examined set

phrases usage was studied diachronically on the base of the big Google books Ngram Viewer Russian corpus counting billions of tokens. The study argues that diachronic productivity is best evaluated with a studying contexts. Used corpus tools enable to do it. Ultimately, it is shown and maintained that corpus linguistics methods and tools allow to create dictionaries of new type which have to include a larger amount of set phrases and collocations than before.

Key words: set phrases, idioms, collocations, collocation dictionaries, corpus, association measures, concordance, diachronical research

Введение

Один из популярных предметов в языкознании — это устойчивые словосочетания. Они изучаются в разных разделах лингвистики и под различными углами зрения. Классическое название для устойчивого сочетания в русской лингвистической терминологии — это фразеологизм. Различные ученые по-разному интерпретируют понятие фразеологизма и его свойства, и существует множество разных классификаций устойчивых сочетаний. Но если их сопоставить, то мы увидим, что перечень этих свойств и сами классификации нередко похожи и пересекаются между собой.

Однако, несмотря на пристальное внимание лингвистов к фразеологии и связанным темам, можно утверждать, что состояние фразеологии сегодня, на наш взгляд, неудовлетворительно. Фразеологический запас русского языка разбросан по разным лексикографическим изданиям, прежде всего это толковые и фразеологические словари, и ни один словарь не может считаться достаточно полным по охвату фразеологического лексикона. Предположительно, словарь такого полного словаря должен насчитывать несколько сотен тысяч единиц. Также нетрудно убедиться, что статьи существующих фразеологических словарей неполны, они плохо структурированы, никак не привязаны к хронологии.

В наши дни эту ситуацию можно существенно улучшить, прибегнув к помощи корпусов. С появлением больших корпусов текстов, в том числе охватывающих длительный промежуток времени, и с появлением в корпусной лингвистике программно-алгоритмических средств, позволяющих оценивать сочетаемость количественно, сложились все предпосылки для создания большого словаря сочетаемости, основанного на корпусах, с количественной параметризацией внутри.

Нужно отметить, что в корпусной лингвистике сложилась методология более широкого понимания фразеологии, и границы фразеологии здесь значительно расширены (или размыты) за счет новых подходов, общим для которых является понятие «статистической устойчивости». Может быть самой знаменитой цитатой в корпусной лингвистике является высказывание Дж. Р. Фёрса «Вы поймете слово по его окружению» (“You shall know a word by the company it keeps”) [Firth 1957: 179]. Там же и тогда же им было введено вошедшее сегодня в широкий оборот понятие коллокации. Об этом же писал И. А. Мельчук.

«Устойчивость сочетания относительно данного элемента измеряется вероятностью, с которой данный элемент предсказывает совместное появление остальных элементов сочетания (в определенном порядке относительно предсказывающего элемента)» [Мельчук 1960: 73].

В корпусной лингвистике в основе методов вычисления силы синтагматической связи между элементами словосочетаний лежат частотные характеристики и структурно-синтаксические модели, на основе которых по формулам так называемых ассоциативных мер (мер ассоциации) вычисляется коэффициент силы связанности или, по-другому, уникальности данного словосочетания.

Как известно, язык — это динамическая система, и это должно находить и находит отражение в словарях и грамматиках. Однако, может быть, меньше всего этот хронологический аспект изучен применительно к фразеологизмам и другим устойчивым сочетаниям. Одна из причин этого — отсутствие до последнего времени больших исторических (диахронических) корпусов.

Данное исследование преследует цель показать, как можно улучшить фразеологические словари и словари сочетаемости корпусными методами. Мы хотим продемонстрировать, как можно пополнить (обогатить) ту или иную словарную статью, где присутствует описание сочетаемости слов, и как можно учесть динамику употребления устойчивых выражений во времени.

1. Материал и инструмент исследования

В качестве материала мы взяли устойчивые сочетания разных типов. Была поставлена задача — расширить «наполнение» исследуемых сочетаний на основе корпусных данных и изучить их «поведение» на большом корпусе в течение длительного промежутка времени.

В качестве материала и инструмента исследования были использованы Национальный корпус русского языка (НКРЯ) (<http://ruscorpora.ru>), корпусы русских текстов ruTenTen 2011 и ruTenTen 2011 sample системы Sketch Engine (<https://the.sketchengine.co.uk/>), корпус русских текстов Araneum Russicum Maius из семейства псевдопараллельных корпусов Aranea Университета им. А. Коменского в Братиславе (<http://ucts.uniba.sk/>) [Benko 2013], русский корпус системы Google books Ngram Viewer и, соответственно, их программные средства. Объем основного корпуса НКРЯ составляет 230 млн словоупотреблений, ruTenTen 2011 sample и русскоязычный Araneum насчитывают по 1200 млн токенов (около 1000 млн текстоформ), ruTenTen 2011 имеет объем более 18 млрд токенов (14,5 млрд текстоформ).

Самый же большой из них корпус русских книг Google books Ngram Viewer (<https://books.google.com/ngrams>). В настоящее время это наиболее мощный инструмент диахронических исследований. Эта информационная система содержит несколько корпусов размеченных текстов книг на 9 языках. На конец 2012 г. база данных насчитывала более 8 млн книг (текстов), что составляет около 6% всех когда-либо опубликованных печатных книг. Корпус книг на русском языке содержит 591 310 текстов общим объемом более 67 млрд

словоупотреблений. Самые поздние публикации, доступные для пользователей в настоящее время, относятся к 2008 году.

Основной лексической единицей (ЛЕ), с которой работает система, является N-грамма последовательность от одной до пяти словоформ. Причем N-грамма, для того чтобы быть учтенной и обработанной, должна встречаться в корпусе не менее 40 раз. Для каждой заданной ЛЕ для заданного временного интервала строится график, по вертикальной оси которого откладывается относительная частота встречаемости заданной N-граммы в корпусе в данном году, выраженная в процентах. На горизонтальной оси показаны годы, входящие в заданный временной интервал. *Каждая кривая графика маркируется цветом, в конце кривой указывается, какой N-грамме (слову или словосочетанию) она соответствует (рис. 1).*

При построении графиков изменения частоты употребления ЛЕ используется так называемое «сглаживание» (smoothing) При нулевом сглаживании в графике учитывается относительная частота встречаемости N-граммы за каждый год. Однако по-настоящему тенденция в динамике встречаемости слов прослеживается более отчетливо при скользящем усреднении данных. Если значение коэффициента сглаживания равно 3, то это означает, что для некоторого года к числу словоупотреблений искомого слова за этот год прибавляется число словоупотреблений за три предыдущих года и три последующих и полученная сумма делится на семь. Относительное значение этой средней величины будет отражено на вертикальной оси.

В системе нет морфологической нормализации ЛЕ, иначе говоря, поиск лексических единиц (слов или словосочетаний), для которых строится график — это поиск по словоформам. Система предусматривает использование пользовательских тегов для модификации условий построения графиков. И в их числе есть тег *_INF (Inflections)*, который строит кривые для всех форм словоизменительной парадигмы данного слова. Однако данная функция для русского языка работает не всегда корректно.

Имеется тег «подстановочный знак» * (wildcard). Ввод его через пробел после N-граммы или до неё позволяет строить график встречаемости десяти наиболее частотных сочетаний данной N-граммы и слова, следующего за ней или ей предшествующего (рис. 8).

Над кривыми графиков возможны операции: суммирование, вычитание, умножение, деление. Например, *суммирование (сложение) кривых*, при котором поисковые слова вводятся в окно запроса через знак +, позволяет суммировать значения каждой точки по оси ординат двух или более кривых. Это может быть использовано как аналог поиска по лемме; например: *лошадь + лошади + лошадей + ... + лошадях*.

Полезная операция *умножение графиков*, позволяющая умножать на n значения всех точек графика (например, *лемматизация*100*). Данная операция позволяет сделать сопоставимым вид кривых, значения которых отличаются на несколько порядков.

Кроме построения графиков, система предоставляет ссылки к текстам, найденным по запросам, где встретились заданные ЛЕ (рис. 9). Как правило,

это библиографические описания книг и фрагменты текстов с выделением в них заданных N-грамм. В некоторых случаях доступен полный текст книги в графическом формате. Более подробно о сервисе Google books Ngram Viewer см. [Захаров, Масевич 2014].

Похожий инструмент под названием «Графики» с 2012 г. работает и в составе НКРЯ. Функционально он подобен сервису Google books Ngram Viewer. Он показывает хронологическое распределение заданных и найденных лексических единиц (словоформ, словосочетаний) в основном корпусе НКРЯ. Вход в этот сервис возможен как со страницы с результатами поиска по произвольному запросу к основному корпусу (ссылка *Распределение по годам*), так и из главного меню (ссылка *Графики*). Можно задавать временные границы. При сходной идеологии, формулы подсчета относительной частоты в сервисах Национального корпуса и Google Ngram Viewer отличаются. Также немного по-другому рассчитывается сглаживание, например, сглаживание, равное 10, усредняет частоту слова с учетом предшествующих и последующих 5 лет, т. е. для данного года берется средняя величина за 11 лет. Имеется возможность показать таблицы с абсолютными и относительными частотами употреблений за каждый год. Из таблиц по гиперссылкам возможен переход к просмотру примеров из корпуса.

2. Эксперименты

В качестве примеров устойчивых сочетаний для исследования были выбраны сочетания трех типов:

- 1) свободные сочетания с характерными определениями к слову «аплодисменты», выражающими, говоря в терминах теории «Смысл Текст», функцию Magn;
- 2) идиома «ничтоже сумняшеся»;
- 3) фразеологизированные сочетания с глаголом «перебиваться» в значении «бедствовать».

2.1. «Аплодисменты»

Посмотрим, какие стандартные определения к слову «аплодисменты» зафиксированы в словарях. Новый Большой академический словарь приводит следующие сочетания: *бурные аплодисменты, гром аплодисментов* [БАС 2004]. Словарь сочетаемости слов русского языка дает: *громкие, продолжительные, долго не смолкающие, несмолкаемые, бурные, дружные, одобрительные, горячие, восторженные, сдержанные, скупые, редкие, жидкие аплодисменты* [Словарь сочетаемости 1983].

В системе Sketch Engine имеется инструмент вычисления коллокаций по 7 мерам ассоциации. В одном из режимов мы получили список из 36 прилагательных, из которых 19 можно считать функцией Magn от слова «аплодисменты».

Вот список этих прилагательных, упорядоченный по алфавиту: *бешеный, бурный, дружный, восторженный, всеобщий, горячий, громкий, несмолкающий, громовой, громогласный, долгий, дружный, неистовый, нескончаемый, несмолкаемый, несмолкающий, оглушительный, продолжительный, шумный.*

В словаре сочетаемости слов русского языка таких определений всего 8. Неизвестно, что отражает порядок их следования. Словарь ничего не говорит о частоте их употребления. Эти данные можно получить из корпусов. Например, НКРЯ дает следующие цифры — см. табл. 1 (здесь и далее количественные данные получены в начале февраля 2015 г.):

Таблица 1. Частота сочетаний
«прилагательное + аплодисменты» в НКРЯ

№ п/п	Словосочетание	Частота в корпусе
1.	бурные аплодисменты	337
2.	продолжительные аплодисменты	125
3.	дружные аплодисменты	81
4.	шумные аплодисменты	59
5.	громкие аплодисменты	47
6.	оглушительные аплодисменты	25
7.	восторженные аплодисменты	22
8.	несмолкающие аплодисменты	2

При этом важно понимать, что и как мы ищем. Так, при поиске в НКРЯ «бурные аплодисменты» находятся в 279 контекстах (интервал 1 слово вправо от *бурные*), в то время как поиск в интервале 3 слова вправо дает нам 337 контекстов. Подавляющая часть прироста обеспечивается сочетанием «бурные и продолжительные аплодисменты». Поиск же по сочетанию «дружные аплодисменты» в интервале 3 слова вправо выдает нам дополнительно не совсем корректные сочетания «дружный смех и аплодисменты», «дружный хохот и аплодисменты» и др. То есть, полученные цифры не нужно абсолютизировать, важны их относительные величины.

Иногда, задав нестандартный режим поиска, можно получить дополнительно интересные результаты. Например, ни один из наших корпусов не дал к *аплодисментам* коллоката *жесткий*. Однако поиск в интервале с отключенным согласованием между этими словами дает фразу «*По жесткому звуку аплодисментов чувствовалось...*», где появляется это определение.

Выявление коллокаций по формулам мер ассоциации позволяет учитывать не только частоту совместной встречаемости, но и частотность или редкость каждого элемента, то есть вычислять силу синтагматической связи между элементами словосочетания. Результат можно упорядочить или по частоте, или по значению меры ассоциации (табл. 2). Система Sketch Engine при этом считает силу связи по разным мерам ассоциации, с учетом структурно-синтаксических формул и с учетом разрывности словосочетаний. И мы видим, что *продолжительные аплодисменты* по силе связи (мера *salience*) оказались лишь на 7 месте.

Таблица 2. Сочетания «прилагательное + аплодисменты» в корпусе ruTenTen 2011, упорядоченные по мере salience

№ п/п	Словосочетание	Частота в корпусе	Мера salience
1.	бурные аплодисменты	13 372	10,25
2.	дружные аплодисменты	2 051	8,42
3.	оглушительные аплодисменты	656	8,29
4.	громкие аплодисменты	3 711	8,22
5.	восторженные аплодисменты	899	8,17
6.	одобрительные аплодисменты	388	8,05
7.	продолжительные аплодисменты	2 495	7,80
8.	несмолкающие аплодисменты	211	7,62

Результаты поиска сочетаний со словом «аплодисменты», полученные разработчиками Генерального Интернет-корпуса русского языка (ГИКРЯ) на их корпусе (подкорпус «Журнальный зал», объем 313 млн словоупотреблений) и любезно предоставленные автору, дают несколько другую картину, а именно: *бурные* 194, *продолжительные* 72, в т.ч. *бурные (и) продолжительные* 33, *громкие* 39, *дружные* 26, *шумные* 17, *восторженные* 15, *долгие* 14, *горячие* 13, *оглушительные* 12, *несмолкающие* 12. Это еще раз говорит о том, что цифры не нужно абсолютизировать и что нужно учитывать, на каком корпусе получены те или иные данные. В данном случае мы имеем, в основном, те же сочетания, но в другом порядке. Это мы уже могли заметить и в табл. 1 и табл. 2, сравнив ранги сочетаний в НКРЯ и в ruTenTen 2011. Однако, если нам важно не сравнение корпусов между собой по тем или иным ЛЕ, а распространенность этих единиц в языке, которую мы хотим определить по частотам в корпусе, то нужно опираться не на абсолютные частоты, а на относительные (ipm). Проиллюстрируем это на маленьком примере (табл. 3).

Таблица 3. Сравнение частот сочетаний в разных корпусах

Словосочетание	Частота в корпусе			ipm		
	НКРЯ	ruTenTen	ГИКРЯ	НКРЯ	ruTenTen	ГИКРЯ
дружные аплодисменты	81	2 051	26	0,350	0,140	0,080
громкие аплодисменты	47	3 711	39	0,200	0,260	0,120
несмолкающие аплодисменты	1	211	12	0,004	0,015	0,038

Как мы видим, разные корпуса по-разному оценивают вес соответствующих сочетаний в языке, а фактически, в подязыке, который представлен каждым корпусом. Отдельная задача — попытаться эту разницу понять и объяснить, с тем

чтобы какие-то особенности корпуса (преобладание какой-то тематики или типа текстов, возможное наличие дублетов и т. п.) не переносить на язык в целом.

Тем не менее, создавая словарь устойчивых сочетаний на основе корпусов, мы имеем возможность выстроить их по частоте употребления или по силе «спаянности», оговорив, на каком материале этот словарь создается. Более того, по-видимому, иногда полезно опираться на усредненные характеристики, полученные на разных корпусах.

Данные, полученные на синхронных корпусах, ничего не говорят о продуктивности этих сочетаний в разные промежутки времени. Чтобы увидеть использование этих сочетаний на протяжении длительного периода, построим графики распределения частоты употребления этих сочетаний в сервисах Google books Ngram Viewer и «Графики» НКРЯ в текстах двух последних столетий. Вот общая картина (рис. 1, рис. 2).

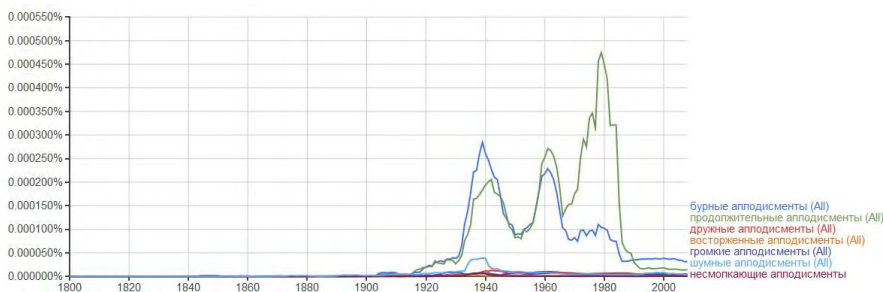


Рис. 1. Кривые встречаемости биграмм со вторым словом «аплодисменты» в корпусе Google books Ngram Viewer

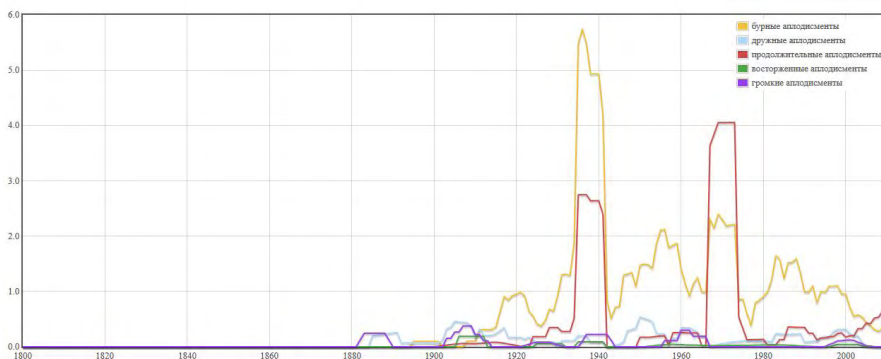


Рис. 2. Кривые встречаемости биграмм со вторым словом «аплодисменты» в корпусе НКРЯ

Как мы уже видели из таблиц 1 и 2, упорядочение по частоте показывает, что «лидируют» *бурные аплодисменты*, а на втором месте идут *продолжительные*. Но это суммарные данные по всему корпусу. На графиках же мы видим распределение частот употребления этих сочетаний во времени см. пики на рубеже 1940-х и 1980-х годов. При этом следует отметить, что в широком употреблении эти сочетания вошли только в XX веке. И если во второй половине 1930-х «верх берут» «бурные», то в конце 70-х — начале 80-х преобладают «продолжительные». Это видно и из сервиса Google, и сервиса НКРЯ. Однако сервис НКРЯ в этом и в других случаях дает картину мало репрезентативную по причине недостаточности данных. Анализируемые словосочетания представлены в корпусе в малых количествах и не в каждом году. Так, *громкие аплодисменты* встретились в основном корпусе НКРЯ по одному разу в 1885, 1906, 1908, 1910, 1925, 1939–40, 1959, 1963, 1998–2000, 2003 гг. и два — в 2001 г. Этого явно мало. Поэтому в дальнейшем мы будем опираться, в основном, на графики системы Google.

Если же задать сглаживание, равное нулю, то можно определить пик использования того или другого сочетания в каждом году (точнее, в текстах данного года; напомним, в корпусе Google books Ngram Viewer это книги) (рис 3).

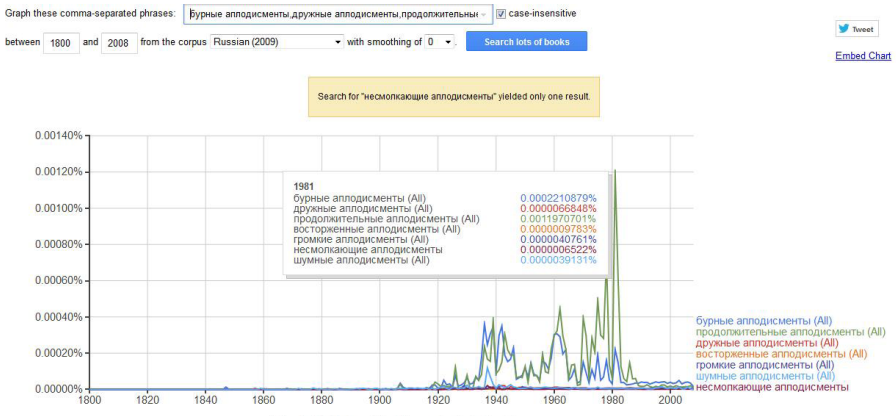


Рис. 3. Кривые встречаемости биграмм со вторым словом «аплодисменты» в корпусе Google books Ngram Viewer со сглаживанием, равным нулю

Проанализируем также атрибутивное отношение, выраженное в форме «существительное в им. пад.+ аплодисменты в род. пад.». Все словари согласно приводят следующие коллокации для выражений этого типа: *шквал*, *гром*, *буря*, *грохот*, *взрыв*. Остается, однако, неясной частотность их употребления. Данные поиска в корпусах ruTenTen 2011 sample и Araneum Russicum Maius и графики показывают явное преимущество коллоката «гром» (рис. 4, рис. 5).

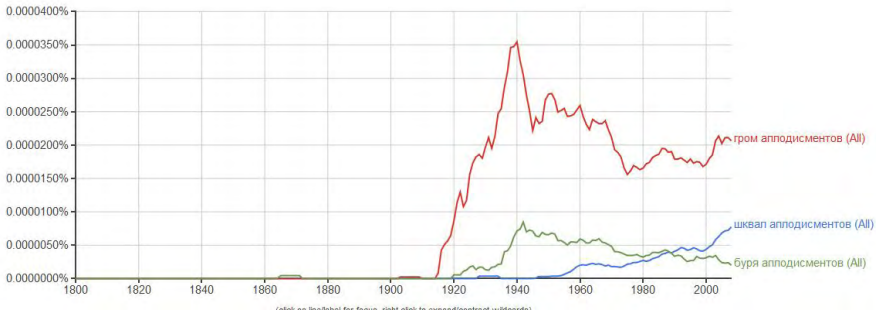


Рис. 4. Кривые встречаемости биграмм с существительным и со вторым словом «аплодисменты» в родительном падеже в корпусе Google books Ngram Viewer

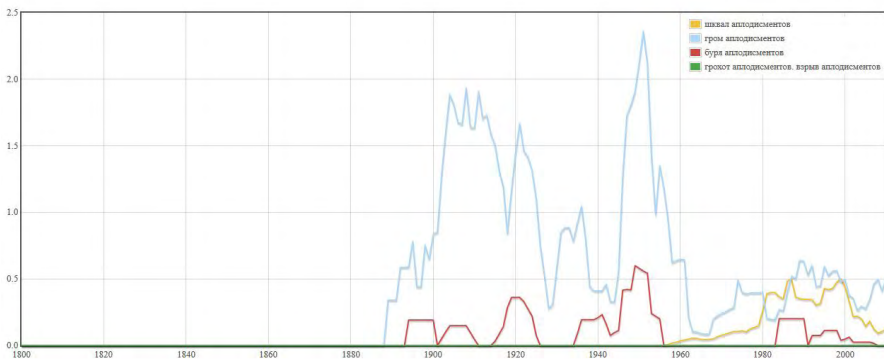


Рис. 5. Кривые встречаемости биграмм с существительным и со вторым словом «аплодисменты» в родительном падеже в корпусе НКРЯ

На втором месте — «буря». И гораздо реже встречаются «грохот» и «взрыв». Но мы видим, что на обоих рисунках, начиная с некоторого периода (1960-е годы по данным НКРЯ и 1980-е по данным Google), «шквал аплодисментов» идет вверх и устойчиво обгоняет «бурю». Но если обратиться к корпусу ГИКРЯ, отражающему современное состояние языка, то в подкорпусах «Живой журнал» и «В контакте» *шквал* уже обошел и *гром*, что говорит о том, что «живой язык» предпочитает последнее сочетание.

На обоих графиках мы видим всплеск в употреблении сочетания «гром аплодисментов» примерно с начала 1900-х по 1930-е годы, и, если обратиться к текстам этого периода, представленным в корпусах, то можно увидеть, что это связано, вероятно, с популярностью театрального искусства, нашедшей отражение в литературе. На рис. 5 второй, еще больший всплеск наблюдается на графике НКРЯ в 1950-е годы, и это связано, по-видимому, с политическими событиями и их отражением в публицистике, богато представленной в НКРЯ.

Поиск в основном корпусе НКРЯ дает 118 вхождений для «гром аплодисментов» (165 с учетом словоизменения, но следует помнить, что графический

сервис строит графики по словоформам), 33 вхождения для «шквал аплодисментов», 15 — для «буря аплодисментов», только одно — для «взрыв аплодисментов» и ни одного для «грохота». На графике Google последних двух сочетаний нет совсем, потому что — напомним — по алгоритму системы каждая N-грамма, чтобы быть представленной на графике, должна встретиться в корпусе не менее 40 раз. Кажется, что *взрыв*, указанный в словаре сочетаемости, в корпусных источниках практически отсутствует.

Однако можем ли мы полностью доверяться корпусным источникам? В 230-миллионном НКРЯ искомые сочетания не нашлись, зато в сегменте «литературного Интернета» (ГИКРЯ, «Журнальный зал») в соизмеримом с НКРЯ по объему корпусе нашлось 25 *взрывов* и 10 *грохотов*. Этот и подобные факты требуют своего объяснения, чтобы мы могли опираться на получаемые результаты либо иногда их отбрасывать. Например, слово или словосочетание с большой частотой может иметь источником «взрывообразное» появление в нескольких текстах в коротком промежутке времени и не быть характерным для языка в целом. Для минимизации таких «всплесков» в лингвистике существуют специальные меры, учитывающие равномерность появления слова в корпусе (коэффициент Жуйана, Average Reduced Frequency и др.).

И последнее. В корпусе *Aganeum Russicum Maius* корпусный менеджер насчитал 14 *взрывов* аплодисментов, в то время как сочетание *взорваться* аплодисментами в разных формах глагола встретилось 42 раза, что говорит о том, что набор конструкций, выражающих функцию Magn от слова «аплодисменты», должен быть расширен.

2.2. «Ничтоже сумняшеся»

Возьмём другой тип устойчивых сочетаний — идиомы и, как пример, сочетание «ничтоже сумняшеся». Словари приводят его в двух формах «ничтоже сумняшеся» и «ничтоже сумняся». Историкам языка это выражение хорошо известно. Но как оно «жило» в языке на протяжении веков, об этом нам расскажет корпус (рис. 6).

И неслучайно у Чехова мы встречаем именно эту форму. «Для нее ясна была эта красивая смелость современного человека, с какою он, не задумываясь и *ничтоже сумняся*, решает большие вопросы и строит окончательные выводы» (А. П. Чехов. *Несчастье*). «Во время своего путешествия из Сахалина я достаточно привык к туманам и свежим ветрам и потому смотрю теперь на Чёрное море свысока и во время качки обедаю *ничтоже сумняся*» (А. Чехов. *Остров Сахалин*). И только во 2-м десятилетии 20-го века стало преобладать «ничтоже сумняшеся».

Мы видим, что на протяжении достаточно долгого времени основной формой было «ничтоже сумняся», так в 1889 г. эта форма в литературе использовалась более чем в 4 раза чаще (рис. 7).

Поиск в сервисе Google позволяет выявить и другие сочетания со словом «ничтоже» (рис. 8), что также представляет интерес для лингвистов, и проанализировать их, перейдя по гиперссылкам собственно к текстам (рис. 9).



Рис. 6. Кривые встречаемости биграммы «ничтоже сумняшся» в корпусе Google books Ngram Viewer

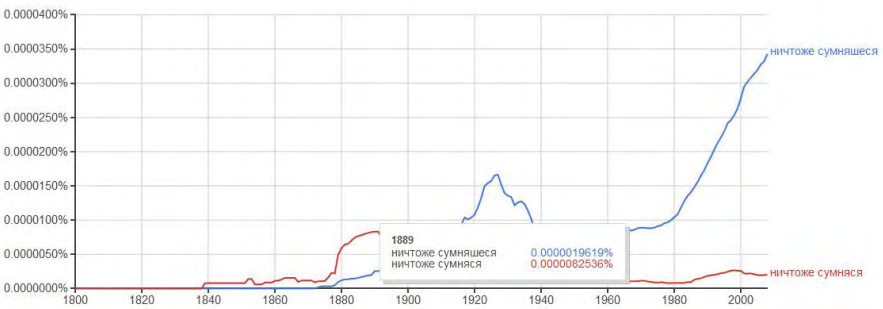


Рис. 7. Кривые встречаемости биграммы «ничтоже сумняшся» в корпусе Google books Ngram Viewer (с относительной частотой встречаемости)

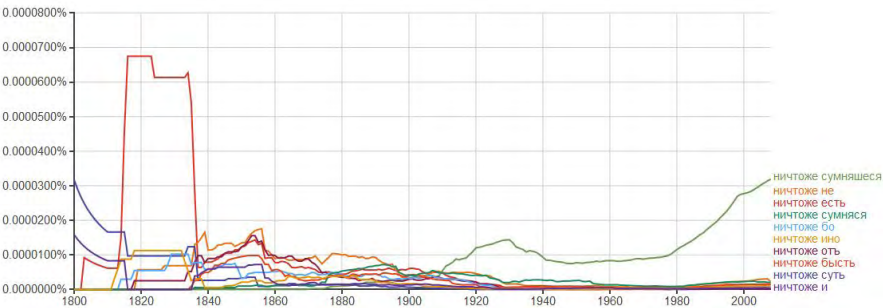


Рис. 8. Кривые встречаемости десяти биграмм с первым словом «ничтоже» (использование подстановочного знака после N-граммы)

[ничтоже есть](#)
[ничтоже сумняшеся](#)
[ничтоже не](#)
[ничтоже суть](#)
[ничтоже оть](#)
[ничтоже бо](#)
[ничтоже бысть](#)
[ничтоже сумняся](#)
[ничтоже ино](#)
[ничтоже и](#)

Рис. 9. Гиперссылки перехода к текстам биграмм с первым словом «ничтоже»

2.3. «Перебиваться»

Многие фразеологизмы и устойчивые сочетания имеют лексико-синтаксические варианты, когда либо меняется лексическое наполнение в рамках некоторой структурной формулы, либо при том же наполнении меняется формула. Например, «кошки скребут». Но где? Словари сообщают, что «на душе» и «на сердце». А где чаще? Выясняется, что чаще «кошки скребут на душе» и больше всего они скребли в лихие годы на переломе 1980–1990-х гг. (рис. 10).

Наверное, не будет ошибкой утверждение, что фразеологизмов с лексико-синтаксическими вариациями большинство. Примеры их можно множить и множить: *беречь (хранить) как зеницу ока; беречь пуще глаза; мерить одной мерой (меркой), мерить на одну меру (мерку); ест за троих, есть в три горла; драть (сдирать/содрать) шкуру (три, две шкуры); драть (сдирать/содрать) по три (две) шкуры; хоть в землю заройся, хоть из-под земли достань; брать/взять (забирать/забрать) в [свои] руки, прибирать/прибрать к рукам; сталкивать/столкнуться лицом к лицу, носом к носу, нос в нос, лоб в лоб* [Бирих, Мокиенко, Степанова 1997]. Такие вариативные сочетания в словарях описаны, естественно, менее полно по сравнению с лексикализованными фраземами.

Рассмотрим этот тип вариаций более подробно на примере сочетания глагола «перебиваться» с предложной конструкцией «с ... на ...». Традиционно словари дают сочетания «*перебиваться с хлеба на квас*», «*перебиваться с хлеба на воду*». Кроме того, приводятся синонимические конструкции «с *куска на кусок*», «с *гроша на копейку*», «с *пуговики на петельку*» [Бирих, Мокиенко, Степанова 1997: 15].



Рис. 10. Кривые встречаемости сочетаний с выражением «кошки скребут» в корпусе Google books Ngram Viewer

Посмотрим, что нам дополнительно дают корпуса, так сказать, качественно и количественно. Поиск контекстов и коллокаций в указанных корпусах добавляет к вышеприведенному списку еще немалую толику, а именно: *с хлеба на воду, с хлеба на кофе, с гроша на грош, с копейки на копейку, с рубля на рубль, с хлеба на квас, с воды на квас, с воды на хлеб, с хлеба на картошку, с петельки на пуговку, со дня на день, с весны на весну, с работы на работу*. Весьма часто встречаются сочетания с «двойки на тройку», а также с «тройки на четверку», но уже, видимо, другое значение. Есть и более экзотические: «с селедки на вермишель», «с „Российского“ на „Докторскую“» (сыр и колбаса). Встречается и такое: «Те, кто в советское время выживал на скудном пайке из грубоватых братьев Вайнеров и предельно идеологизированного Юлиана Семенова, а позже несколько лет *перебывался с умеренно культурной Дашковой на пещерную Серову*, радостно влились в ряды поклонников Бориса Акунина и Леонида Юзефовича». «Раз уж мы в беседе постоянно *перебываемся с воды на снег*, ответьте по удобному случаю, почему серебряный призер чемпионата России в трековых гонках не сумел победить в традиционном мартовском каракулинском эндуро?» А к вершинам народного языкового творчества можно отнести вот это: «И это беспроводной интернет! МТС перешагнул порог в 1 Мбит/с — 1164 Кбит/с. У Мегафона четкий прием — 3894 Кбит/с. Мне, молившемуся на dial-up, который *перебывался с 22 на 40 килобит в секунду*, это кажется чем-то фантастическим».

Вышеприведенные примеры в основной массе получены из корпусов, созданных по технологии Wasky, то есть на основе текстов из веба, и нередко они отражают языковое творчество, но не узус, подчеркивая лишь продуктивность конструкции «перебываться с ... на ...». Однако наличие больших корпусов позволяет выявлять кандидатов на вхождение в словарь, а статистические данные, получаемые на них, дают возможность оценить широту и сферу распространенности того или другого сочетания.

Наиболее частые сочетания для данной конструкции в корпусе Google books Ngram Viewer показаны на рис. 11.



Рис. 11. Кривые встречаемости выражений с глаголом «перебываться» в корпусе Google books Ngram Viewer

Из графика видно, что наиболее частые устойчивые сочетания с глаголом «перебываться» это те, которые приводятся в фразеологических словарях, и что активно в языковой обиход они вошли только в 20-м веке. А вышеприведенные примеры в основной массе получены из корпусов, созданных по технологии Wasky, то есть на основе текстов из веба, и отражают языковое творчество, но не узус.

И здесь еще раз следует подчеркнуть, что, интерпретируя корпусные данные, мы должны хорошо понимать, что собой представляет тот или другой корпус и как эти данные получаются. Например, данные анализа по корпусу ГИКРЯ, предоставленные автору его разработчиками, показывают, что если в корпусе книг Google сочетание с хлеба на квас встречается чаще, чем с хлеба на воду, то во всех трех подкорпусах ГИКРЯ картина диаметрально противоположная: в сумме 269 употреблений *с хлеба на воду* против 97 *с хлеба на квас*. То же соотношение демонстрирует корпус ruTenTen 2011 (787 против 370). Все это позволяет говорить о различии в использовании этих выражений между книжным языком и современным «спонтанным».

3. Опыт корпусного исследования

Проведем небольшое исследование, как выглядит в корпусах сочетаемость слова «мастер». Сочетаемость слов определяется различными факторами: лексическими, грамматическими, семантическими, стилистическими. Все они влияют на норму и на узус. Можно утверждать, что узус — один из определяющих факторов при составлении словарей включать какую-либо единицу в словаре или нет, какие-примеры подобрать в качестве речений,

какие имеются для данного слова идиомы. Один из подходов в изучении узуса, получивший распространение в последние годы, заключается в выявлении статистических закономерностей на корпусах текстов.

Слово «мастер» в русском языке является довольно частотным. Его *ipm* по электронной версии Частотного словаря современного русского языка (на материалах Национального корпуса русского языка) О. Н. Ляшевой и С. А. Шарова (<http://dict.ruslang.ru/freq.php>) равняется 100,8, в корпусе *ruTenTen 2011* он равен 96,0. В разных словарях *мастер* обычно имеет 4–7 значений и его сочетаемостные свойства никак особенно не описываются.

Рассмотрим его «поведение» на материале корпусов, которые были описаны выше. При изучении конкордансов с этим словом обращаешь внимание на сочетание «каких-то дел мастер». И таких контекстов, кроме привычного «золотых дел мастер», в корпусах находится достаточно много (329 в НКРЯ и 4015 в *ruTenTen 2011*). Если объединить все определения к словоформе «дел», то их будет 252:

абордажных, автомобильных, аккордеонных, алмазных, багетных, баланных, банкетных, банных, барабанных, баррикадных, бархатного, берестяных, библиотечных, бриллиантовых, броневых, бронзовых, бронзовых, бронно-кольчужных, булочных, бумажных, буровых, бытовых, взрывных, винных, витражных, водочных, воровских, выборных, вывесочных, газетных, газовых, гаражных, гармонных, гитарных, глазных, гламурных, глиняных, гончарных, горных, городовых, грильных, гробовых, дамских, дверных, деревянных, деспотических, дипломатических, добрых, домашних, дорожных, железных, жестяных, живописного, журнальных, закулисных, замочных, заплатных, заплечных, запретных, здоровых, земельных, зеркальных, золотых, игрушечных, изразцовых, именных, иностранных, искусных, кабельных, кальянных, каменных, каминных, каретных, карманных, картежных, кирпичных, ключных, книжных, ковровых, кожаных, кожевенных, колбасных, колесных, колодезных, колокольных, колыбельных, кольчужных, комедийных, компьютерных, конкурсных, конфетных, коньячных, копировальных, корабельных, корабельных, костяных, кофейных, красочных, крепостных, крепостных, кроватных, кровельных, кровопийственных, кузнечных, кузовных, кукольных, кулачных, кулинарных, кухонных, ледяных, лепных, литейных, литературных, лодочных, любовных, макетных, малярно-живописных, малярных, машинных, мебельных, мебельных, медных, мироедских, мозаичного, мозольных, молочных, монетных, мостовых, музыкальных, музыкальных, мусийных, мясных, ножевых, обувных, огненных, оконных, оловянных, оптических, органичных, оружейных, открыточных, палатных, палаточных, палаческих, памятных, парусных, переговорных, переплетных, персонных, перспективных, перчаточных, песочных, печатных, печных, плотницких, погребальных, поддельных, подкопных, подручных, подъемных, позолотных, половых, помойных, портновских, портных, портняжных, похоронных, почтовых, преоспективного, прикладных, пробирных, прохвостных, пушечных, пыточных, пытошных, ракетных, резных, рекламных, ресторанных,

ритуальных, розыскных, ручных, рыбных, садовых, самолетных, сапожных, сателлитных, седельных, селедочных, сердечных, серебряных, сих, скрипичных, сладких, слесарных, словесных, социальных, ссудных, стегательных, стеклянных, стекольных, стекольных, столярных, страховых, строительных, суконных, сусалнаго, сценических, сыскных, табачных, табуреточных, тайных, телевизионных, ткацких, токарных, топорных, трубных, угольных, ударных, фальшивых, фейерверкского, фершельных, фискальных, фонтанных, фотографических, фотошопных, хлебных, ходульных, холодильных, хрустальных, цветочных, ценных, цеховых, циркульных, чайных, часовых, чеканного, чемоданных, черепаховых, чернильных, шапочных, шашлычных, швейных, шлифовальных, шлюзных, шляпных, шляпочных, шоколадных, ювелирных, янтарных.

Анализ показывает, что за исключением немногих, окказиональных или оценочных, все они относятся к какому-либо ремеслу. Посмотрим на частотные характеристики этих сочетаний. Приведем 25 наиболее частотных сочетаний, полученных на корпусе ruTenTen (табл. 4).

Таблица 4. Наиболее частотные прилагательные в сочетаниях типа «таких-то дел мастер», упорядоченные по частоте совместной встречаемости

Ранг	Слово	Частота
1.	золотых	519
2.	часовых	261
3.	серебряных	179
4.	заплечных	112
5.	каменных	111
6.	кукольных	92
7.	гробовых	73
8.	пыточных	58
9.	витражных	57
10.	добрых	49
11.	оружейных	46
12.	кузнечных	46
13.	кузнечных	46

Ранг	Слово	Частота
14.	колокольных	44
15.	ювелирных	40
16.	шляпных	39
17.	чемоданных	38
18.	обувных	38
19.	похоронных	37
20.	мебельных	37
21.	деревянных	36
22.	сапожных	33
23.	печатных	31
24.	скрипичных	30
25.	столярных	27

Высокая величина частоты совместной встречаемости, казалось бы, говорит об устойчивости данного сочетания. Однако этой характеристики недостаточно, чтобы говорить о предпочтительной сочетаемости одного слова с другим. Вполне возможно, что при невысокой частоте совместной встречаемости сочетание представляет собой неделимое единство. Имеется целый ряд статистических мер, получивших название меры ассоциации, или меры ассоциативной связанности (англ. association measures), вычисляющих

силу связанности элементов в составе сочетаний. Значения мер ассоциации можно считать показателями силы синтагматической связи между элементами словосочетаний. Приведем те же 25 сочетаний с подсчитанными значениями нескольких мер ассоциации (табл. 5).

Таблица 5. Наиболее частотные прилагательные в сочетаниях типа «таких-то дел мастер», упорядоченные по значению меры ассоциации MI3

Ранг по частоте	Слово	Частота сочетания	MI3	log likelihood	logDice	MI.log_f	Ранг по MI3
9.	витражных	57	33,096	1472,831	8,651	85,745	1.
1.	Золотых	519	32,046	9180,223	6,822	87,557	2.
4.	Заплечных	112	31,295	2132,625	9,131	82,657	3.
7.	Гробовых	73	30,744	1421,884	8,719	77,587	4.
2.	Часовых	261	30,188	4034,095	7,368	78,951	5.
17.	чемоданных	38	29,400	822,664	7,993	68,098	6.
13.	кузнечных	46	28,675	520,601	7,321	60,985	7.
20.	мебельных	37	28,365	594,114	7,550	61,788	8.
16.	шапных	39	28,330	771,052	7,909	64,324	9.
3.	серебряных	179	28,301	2676,989	6,488	69,416	10.
8.	пыточных	58	28,261	971,907	8,109	66,014	11.
6.	кукольных	92	27,395	1239,558	7,604	64,266	12.
14.	колокольных	44	27,085	799,089	7,754	60,775	13.
23.	печатных	31	27,023	584,485	7,532	58,078	14.
22.	сапожных	33	26,760	556,710	7,464	56,926	15.
24.	скрипичных	30	25,642	534,014	7,264	53,529	16.
5.	каменных	111	25,584	1529,471	5,062	56,727	17.
10.	добрых	49	25,089	650,132	6,918	53,147	18.
11.	оружейных	46	25,050	719,464	6,724	53,500	19.
21.	деревянных	36	25,019	414,504	7,010	49,856	20.
12.	кузнечных	46	24,712	500,269	6,942	50,463	21.
19.	похоронных	37	24,203	510,628	6,642	49,124	22.
18.	обувных	38	23,816	509,175	6,364	47,956	23.
25.	столярных	27	22,038	344,005	5,640	40,968	24.
15.	ювелирных	40	20,973	423,087	3,777	38,368	25.

Мы видим, что разные меры по-разному оценивают силу синтагматической связи и что ранги сочетаний в списках, упорядоченных по значениям мер, не обязательно совпадают с рангами в списке, упорядоченном по частоте совместной встречаемости. И не всегда большую силу связи получают наиболее частые сочетания, например, сочетание «витражных дел мастер», всего лишь девятое по частоте (см. выше табл. 4), оказывается первым по рангу меры MI3

и, по-видимому, с бóльшим основанием может быть включено в словарь в качестве (или как пример) устойчивого сочетания. В корпусной лингвистике устойчивые сочетания могут быть определены как *статистически* устойчивые словосочетания. При этом статистически устойчивое сочетание может быть как фразеологизированным, так и свободным.

В лексикографических изданиях следует, вероятно, указать, что слово «дело» в этих сочетаниях почти всегда стоит во множественном числе (из 329 сочетаний в НКРЯ только в 17 *дело* в единственном числе). Также нужно упомянуть, что в этого рода сочетаниях имеет место именно такой порядок слов. Конечно, не будет ошибкой сказать, «мастер кузнечных дел», но так не говорят. Или все-таки говорят? Только проверка на большом «живом» текстовом материале даст нам ответ. Вот что показывает корпус ruTenTen 2011: «мастер золотых дел» 12 сочетаний против 579 с «мастером» в постпозиции, «мастер серебряных дел» 12 против 179, «мастер заплочных дел» 15 и, соответственно, 112, «мастер гробовых дел» 2 и 73, «мастер пыточных дел» 1 и 58.

Стоит обратить внимание на сочетание «сих дел мастер». В НКРЯ оно встречается 4 раза, в корпусе ruTenTen 2011 12 раз, но и во всем вебе в Яндексе — всего 722 раза, причем это все дубли, а разных цитат немногих более 12. Наиболее часто представлено в вебе высказывание Л. Троцкого о В. Ленине из письма к Н. Чхеидзе, где он пишет о склоке, «которую разжигает *сих дел мастер* Ленин, этот профессиональный эксплуататор всякой отсталости в русском рабочем движении». Изредка это сочетание встречается в литературе, в частности, у Н. К. Михайловского, А. П. Чехова, Н. А. Тэффи, С. В. Максимова, Н. Е. Врангеля, И. Н. Потапенко, В. Ф. Пановой, причем нередко в кавычках. На самом деле это сочетание того же профессионально-ремесленного происхождения, что и вышеприведенные выражения. Вот как это выглядело в жизни, читаем: «Как в Киеве я смеялся, смотря на вывеску, на которой были изображены самовар, мельница и ножницы и подписано: «сих дел мастер»...».

Также для анализа сочетаемости представляют интерес предложные сочетания со словом «мастер» (прежде всего с предлогами «по», «на», «с», «в», «от»), где оно выступает как «хозяин» в структуре зависимостей, то есть управляет предлогом, который связывает его со знаменательным словом.

Рассмотрим здесь только сочетания с предлогом «на». Очень часто это синтаксема (отношение между хозяином, предлогом и слугой), которую Г. А. Золотова определяет как потенсив — «синтаксема от отвлеченных имен, обозначающих потенциальное действие при словах модальной семантики (глаголах, именах, прилагательных). С личными именами (мастер, мастак, охотник) потенсив образует сочетание, представляющее собой модальную и экспрессивно-оценочную модификацию предикативной характеристики лица» [Золотова 1988: 197]. Отметим, что для таких сочетаний имеется синонимичная конструкция с глаголом (*мастер на шутки* — *мастер шутить*).

Какие же «процессные существительные» встречаются в корпусах в сочетании с *мастером*? В корпусе ruTenTen 2011 находится 34 сочетания с предлогом «на» и со словом «дело», причем «дело» всегда стоит во множественном числе и всегда с определением (за исключением одного случая). Частые

определения к делам *эти* и *такие*, кроме того, встретились *темные*, *плохие*, *маленькие* и *пытошные*.

В числе других «потенциальных действий» для слова «мастер» были найдены: *штуки*, *штучки*, *шутки*, *проделки*, *операции*, *авантюры*, *интриги*, *трюки*, *разговоры*, *анекдоты*, *выдумки*, *флешмобы*, *проказы*, *хитрости*, *слова*. Почти всегда с определениями, среди которых преобладают *такие*, *всякие*, *подобные*, *разные*, *всевозможные*, также встречаются *сходственные*, *веселые*, *подлые*, *хаккерские*, *жестокие*, *недобрые*. В 15 случаях следом за предлогом идут сочетания *всякого* рода, *такого* рода, *разного* рода, *различного* рода.

Встречаются в корпусах и фразеологизированные выражения со словом «мастер»:

- *мастер на все руки* (78 раз в НКРЯ, 2910 раз в ruTenTen 2011);
- *мастер с большой буквы* (4 вхождения в НКРЯ и 470 в ruTenTen 2011, еще 7 вхождений в ruTenTen со следующими определениями к букве: *самая большая*, *высокая*, *огромная*, *маленькая*, *та самая*, *вышшая*);
- *дело мастера боится* (30 вхождений в НКРЯ, 260 в ruTenTen 2011, с определениями *всякое*, *любое*, *ночное дело*).

Интересны сочетания *мастера* с выражениями *на свой лад*, *вкус*, *глаз* — в этом случае *мастеру* всегда предшествуют определения *всякий*, *всяк*, *каждый*.

Есть и другие сочетания для слова *мастер* и с другими предлогами, но на этом мы здесь остановимся.

4. Заключение и выводы

Сегодня русский язык переживает период быстрого обновления своего состава. Не избежали этого и устойчивые сочетания. На периферии языка оказываются сочетания, отражающие некоторые стороны социальной жизни до-революционного общества (мир чиновничества, картежные игры и др.). Зато повышенную частотность получают сочетания из области науки, техники, спорта. Чтобы увидеть все эти изменения, нужны большие корпуса, особенно для фразеологии, учитывая сравнительно низкую частоту употребления фразеологизмов в текстах. И сейчас такие корпуса начинают появляться.

В то же время необходимо, чтобы корпусная лингвистика развивала свои средства. Так, система Google Books Ngram Viewer предоставляет большие возможности для историко-культурных и лингвистических исследований. Однако в текстах корпуса встречается много ошибок распознавания. Поиск заданных лексических единиц ведется по словоформам, а не по леммам. Корпус построен исключительно на книгах и тем самым не сбалансирован. По-видимому, целесообразно было бы провести основательное исследование с применением методов статистической обработки данных, чтобы понять, как эти и другие проблемы влияют на достоверность получаемых результатов. Все это относится и к сервису НКРЯ «Графики» (главный недостаток малый для полноценных диахронических исследований объем корпуса).

Проведенное исследование показало, что корпуса и инструментарий корпусной лингвистики позволяют выявить и существенно расширить лексический фонд устойчивых словосочетаний разного типа и особенности их бытования. Основываясь на корпусах, лингвисты имеют возможность создавать словари и учебники нового типа, где сочетаемость будет представлена неизмеримо шире, чем до сих пор. В качестве примера такого словаря можно привести словарь «КроссЛексика», в котором словосочетания составляют самую важную и самую объемную его часть (2,26 млн словосочетаний) [Большаков 2009]. При этом они должны иметь количественные характеристики как силы устойчивости в синхронии, так и историю их употребления в диахронии.

В ходе исследования мы также неоднократно убеждались, что для того, чтобы можно было делать достоверные выводы на основе корпусных данных, следует хорошо представлять себе недостатки и ограничения тех инструментов, которыми мы пользуемся.

Литература

1. БАС — Большой академический словарь русского языка. Том 1. М.—СПб.: Наука, 2004.
2. Беликов В. И., Копылов Н. Ю., Пиперски А. Ч., Селегей В. П., Шаров С. А. Корпус как язык: от масштабируемости к дифференциальной полноте // Компьютерная лингвистика и интеллектуальные технологии: по материалам ежегодной Международной конференции «Диалог» (Бекасово, 29 мая — 2 июня 2013 г.). Вып. 12 (19). М.: Изд-во РГГУ, 2013. С. 84–95.
3. Бирх А. К., Мокиенко В. М., Степанова Л. И. Словарь фразеологических синонимов русского языка. Ростов-на-Дону, 1997.
4. Большаков И. А. КроссЛексика — большой электронный словарь сочетаний и смысловых связей русских слов // Компьютерная лингвистика и интеллектуальные технологии. Международная конференция «Диалог 2009». Вып. 8 (15) М.: Изд-во РГГУ, 2009. С. 45–50.
5. Захаров В. П., Масевич А. Ц. Диахронические исследования на основе корпуса русских текстов Google books Ngram Viewer // Структурная и прикладная лингвистика. Вып. 10. СПб.: Изд-во С.-Петербургского ун-та, 2014. С. 303–327.
6. Золотова Г. А. Синтаксический словарь. М.: Наука, 1988.
7. Мельчук И. А. О терминах «устойчивость» и «идиоматичность» // Вопросы языкознания. 1960, № 4. С. 73–80.
8. Словарь сочетаемости слов русского языка / Институт русского языка им. А. С. Пушкина; Под ред. П. Н. Денисова, В. В. Морковкина. — 2-е изд., испр. — М.: Русский язык, 1983.
9. Benko V. Aranea: Yet Another Family of (Comparable) Web Corpora. In: Petr Sojka, Aleš Horák, Ivan Kopeček and Karel Pala (Eds.): Text, Speech and Dialogue. 17th International Conference, TSD 2014, Brno, Czech Republic, September 8–12, 2014. Proceedings. LNCS 8655. Springer International Publishing Switzerland, 2014, pp. 257–264. ISBN: 978-3-319-10815-5.

10. *Firth, J. R.* 1957. A synopsis of linguistic theory 1930–1955. In: F. Palmer (Ed.), *Selected Papers of J. R. Firth 1952–1959*. London: Longman, pp. 168–205.

References

1. *BAS* (2004), *Great Academic Dictionary of the Russian Language*, [Bol'shoy akademicheskiy slovar' russkogo yazyka], vol. 1, Moscow/Saint-Petersburg, Nauka.
2. *Belikov V., Kopylov N., Piperski A., Selegey V., Sharoff S.* (2013), *Corpus as language: from scalability to register variation*, [Korpus kak yazyk: ot masshtabiruyemosti k differentsial'noy polnote], *Computational Linguistics and Intellectual Technologies: Proceedings of the International Conference "Dialog 2013"* [Komp'yuternaya lingvistika i intellektual'nye tekhnologii: po materialam ezhegodnoy mezhdunarodnoy konferentsii "Dialog 2013"], vol. 12 (19), Moscow, RGGU, pp. 84–95.
3. *Benko V.* (2014), *Aranea: Yet Another Family of (Comparable) Web Corpora*, In: Petr Sojka, Aleš Horák, Ivan Kopeček and Karel Pala (Eds.): *Text, Speech and Dialogue. 17th International Conference, TSD 2014, Brno, Czech Republic, September 8–12, 2014. Proceedings. LNCS 8655*. Springer International Publishing Switzerland, pp. 257–264, ISBN: 978-3-319-10815-5.
4. *Birikh A. K., Mokiyeiko V. M., Stepanova L. I.* (1997), *Dictionary of phraseological synonyms of the Russian language*, [Slovar' frazeologicheskikh sinonimov russkogo yazyka], Rostov-on-Don.
5. *Bolshakov I. A.* (2009), *CrossLexica: a large electronic dictionary of collocations and semantic links between Russian words*, [KrossLexika — bol'shoy elektronnyy slovar' sochetaniy i smyslovykh svyazey russkikh slov], *Computational Linguistics and Intellectual Technologies: Proceedings of the International Conference "Dialog 2009"* [Komp'yuternaya lingvistika i intellektual'nye tekhnologii: Trudy mezhdunarodnoy konferentsii "Dialog 2009"], vol. 8 (15), Moscow, RGGU, pp. 45–50.
6. *Collocability Dictionary of Russian Language Words* (1983), [Slovar' sochetayemosti slov russkogo yazyka], P. N. Denisov, V. V. Morkovkin (eds.), Moscow, Russkiy yazyk.
7. *Firth, J. R.* (1957), A synopsis of linguistic theory 1930–1955, In: F. Palmer (Ed.), *Selected Papers of J. R. Firth 1952–1959*, London, Longman, pp. 168–205.
8. *Mel'chuk I. A.* (1960), About the terms steadiness and idiomaticity, [O terminakh ,ustoyvchivost' i ,idiomatichnost'], *Questions of Linguistics*, [Voprosy yazykoznaniya], 1960, No. 4, pp. 73–80.
9. *Zakharov V. P., Masevich A. Ts.* (2014), *Diachronic researches on the base of the Russian Google books Ngram Viewer text corpus* [Diakhronicheskiye issledovaniya na osnove korpusa russkikh tekstov Google books Ngram Viewer], *Structural and Applied Linguistics* [Strukturnaya i prikladnaya lingvistika], vol. 10, Saint-Petersburg, pp. 303–327.
10. *Zolotova G. A.* (1988), *Syntactic dictionary* [Sintaksicheskiy slovar'], Moscow, Nauka.