

КОНТРАСТИВНЫЕ ИССЛЕДОВАНИЯ И ПРИКЛАДНАЯ ЛИНГВИСТИКА

Материалы
Международной научной конференции
г. Минск, 29–30 октября 2014 г.

В двух частях
Часть вторая

**ТЕХНОЛОГИИ ИЗВЛЕЧЕНИЯ ИНФОРМАЦИИ ИЗ ИНТЕРНЕТА
И ЕЕ ОБРАБОТКА**

<i>Бочкова А. Л.</i> Автоматическое резюмирование мнений участников интернет-коммуникации об определенном объекте	3
<i>Грибач Т. А.</i> Технология Data Mining для поиска решений	7
<i>Дзубук Т. В.</i> Подходы к автоматическому извлечению фактов из неструктурированного текста	11
<i>Семак А. С.</i> Подходы к автоматическому определению субъективности в интернет-текстах	15
<i>Товпинец К. Г.</i> Основные трудности анализа тональности текстов	21
<i>Чалагаева В. В.</i> Языковая репрезентация категории тональности в новостных статьях англо- и русскоязычной прессы	25
<i>Bilan V., Bobkov A., Gafurov S., Krasnoproshin V., Laar J. van de, Vissia H.</i> <i>An Ontology-based Approach to Opinion Mining</i>	27
<i>Bobkov A., Gafurov S., Krasnoproshin V., Romanchik V., Vissia H.</i> <i>Information Extraction Based on Semantic Patterns</i>	33
<i>Russo Claudio.</i> Extracting linguistic Data from Usenet newsgroups: troubles and challenger	39

АВТОМАТИЧЕСКИЙ АНАЛИЗ И СИНТЕЗ ТЕКСТОВ

<i>Авдеева Н. А., Боярский К. К.</i> Обработка синтаксических групп с числительными в системе автоматического анализа текста SEMSIN.....	43
<i>Анисимович А. В.</i> Задача выделения слова в автоматической обработке текстов на китайском языке	47
<i>Детскина Р. В.</i> Сложные прилагательные в немецком языке и особенности их автоматического перевода на русский язык	50
<i>Елисеева О. Е.</i> Гипертекстовая модель сетевого учебного словаря иностранного языка	55
<i>Зубова И. И.</i> Морфология мифа как основа системы автоматического порождения мифологического текста	59
<i>Кошчанка У. А.</i> Актуальнасць стварэння лексіка-семантычнай анталогії тыпу WordNet для беларускай мовы.....	64
<i>Лысы С. І., Гецэвіч Ю. С.</i> Частотны анализ электроннага тэксту на прадмет выкарыстання слоў і іншых сімвалных паслядоўнасцяў з дапамогай www.corpus.by	68
<i>Метлицкая Н. А.</i> Формализованное представление основного статического содержания текста	72
<i>Нагорнович Э. В.</i> Способы автоматического исправления ошибок в письменном тексте	75

<i>Праблемы амографіі з дапамогай NOOJ для больш чым 50 амографаў рускай мовы</i>	79
<i>Романов Ю. В., Товмач Ю. В.</i> Оценка погрешности определения лингвостатистических параметров	83
<i>Чапля А. И., Чапля С. Г.</i> Опыт статистического разграничения омонимии и полисемии слов на базе триад	87
<i>Черепович А. В.</i> К проблеме использования нормативных словарей при проведении криминалистической экспертизы звукозаписей	88
<i>Шмыга Д. Л.</i> Формальная модель системы автоматической генерации текстов англоязычной контекстной рекламы	92
<i>Шуманская Т. А.</i> Параўнальны аналіз электронных слоўнікаў беларускай мовы	96
<i>Шустова Д. С.</i> Автоматическая экспертиза текстов англоязычных СМИ с элементами вербальной агрессии	100
<i>Яковинин В. С.</i> О методе алгоритмического определения грамматических классов слов	105

КОРПУСНЫЕ ИССЛЕДОВАНИЯ

<i>Беридзе М. М., Бакурадзе Л. Д.</i> Грузинский диалектный остров в Иране	108
<i>Беридзе М. М., Киквидзе З. З., Лорткипанидзе Л. Л.</i> Грузинский корпус метаязыка лингвистики: проблемы и решения	111
<i>Бускунбаева Л. А., Сиразитдинов З. А., Ишмухаметова А. Ш.</i> О внедрении башкирского языка в киберпространство	115
<i>Детскина Р. В., Василевская В. А., Василенко Е. С., Жданович А. Е., Филимонова Т. А.</i> Сопоставительный анализ лексики школьных учебников по белорусскому и русскому языкам	119
<i>Детскина Р. В., Шимчук Л. В.</i> Особенности формирования лексики компьютерного сленга	120
<i>Жубанов А. К., Фазылжанова А.</i> О концепции по созданию Национального корпуса казахского языка (НККЯ)	123
<i>Жубанов А. К.</i> Принципы автоматизации морфологической разметки текстов Национального корпуса казахского языка	126
<i>Зыгмантович Н. В., Ильина Е. И., Марковская Е. В., Шалимо Н. В.</i> Об основных задачах изучения лексики школьных учебников по белорусской и русской литературе	131
<i>Камишлова О. Н.</i> Контрактивные исследования в учебном корпусе текстов	132
<i>Кривошечя И. А.</i> Возможности параллельных корпусов текстов при переводе	135
<i>Лорткипанидзе Л. Л., Еремьян Р. А.</i> Разработка менеджера корпуса грузинских литературных текстов	138

МИНИСТЕРСТВО ОБРАЗОВАНИЯ РЕСПУБЛИКИ БЕЛАРУСЬ
Минский государственный лингвистический университет

КОНТРАСТИВНЫЕ ИССЛЕДОВАНИЯ
И
ПРИКЛАДНАЯ ЛИНГВИСТИКА

МАТЕРИАЛЫ
Международной научной конференции
г. Минск, 29–30 октября 2014 г.

В двух частях
Часть вторая

Минск 2015

УДК 81'42 + 81'33

ББК 81.04 + 81.1

K65

Рекомендованы Редакционным советом Минского государственного лингвистического университета. Протокол № 3 (19) от 23.12.2014 г.

Редакционная коллегия: А. В. Зубов (*отв. редактор*), Т. П. Карпилович (*отв. редактор*), Л. Н. Беляева, М. Г. Богова, Р. В. Детскина, Н. Ю. Павловская, И. В. Совпель

Рецензенты: доктор филологических наук, профессор Д. Г. Богушевич (МГЛУ); кандидат филологических наук, доцент А. И. Головня (БГУ)

K65 Контрастивные исследования и прикладная лингвистика : материалы Междунар. науч. конф., Минск, 29–30 окт. 2014 г. В 2 ч. Ч. 2 / отв. ред.: А. В. Зубов, Т. П. Карпилович. – Минск : МГЛУ, 2015. – 227 с.

ISBN 978-985-460-669-9 (Ч. 2).

ISBN 978-985-460-658-3.

В издание, состоящее из двух частей, включены материалы, ракрывающие возможности контрастивного изучения самых различных языков: белорусского, русского, английского, немецкого, башкирского, туркменского, итальянского, китайского, грузинского, казахского, польского. Показано как результаты такого изучения используются при обучении, автоматической обработке текстов, их переводе.

Данная вторая часть книги содержит тексты выступлений участников конференции в секциях «Технологии извлечения информации из Интернета и ее обработка», «Автоматический анализ и синтез текстов», «Корпусные исследования», «Контрастивная лингвистика и вопросы преподавания языков» и «Особенности речевого поведения в разных культурах».

Для научных работников, аспирантов и магистрантов.

УДК 81'42 + 81'33

ББК 81.04 + 81.1

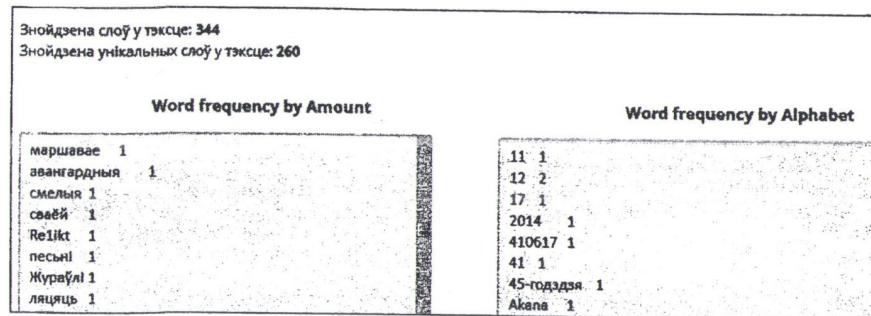
ISBN 978-985-460-669-9 (Ч. 2)
ISBN 978-985-460-658-3

© УО «Минский государственный
лингвистический университет», 2015

Праца з тэкстам, дзе сустракаюцца розныя алфавіты, лічбы і іх камбінацыі, выглядзе так: Адкрываем тэкст па спасылцы <http://news.tut.by/culture/110617.html>. Адраву ж бачым, што ў ім выкарыстоўваюцца як беларускія (дымек, рэалізаваць, праект), рускія (комментарыев), англійскія (Detroit, Naka, Tonqikod) слова, так і лікі (11, 17), змешаныя адзінкі (Re1ikt, 17:41) і іншыя (Re:Песняры). Прыклады настроек сэрвісу, каб атрымаць рапшэнне некаторых камп'ютарна-лінгвістычных задач, могуць выглядаць так.

Задача 1. Знайсці ўсе слова тэксту і падлічыць іх частату выкарыстання.

Рашэнне. Вокны са спісамі сімвалоў не змяняем. Націкаем кнопку «Атрымаць частату слоў!». Бачым рапшэнне на мал. 2.



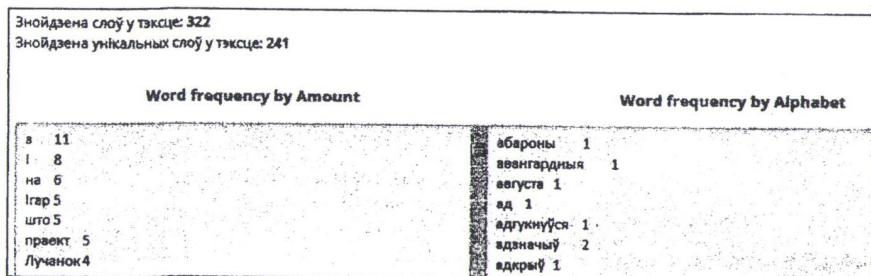
Мал. 2. Рацэнне задачы 1

У атрыманых спісах сімвольных паслядоўнасцяў даследчык можа зauważыць слова (маршавае, складаных, Akana), лікі (11, 2014), камбінацыі слова і лічбаў для стварэння назваў (Re1ikt) ці скарачэнняў (45-годдзя). Прааналізуем іншыя спісы сімвольных паслядоўнасцяў.

Задача 2. Знайсці ўсе слова тэксту, запісаныя кірыліцай, падлічыць частату іх ўжывання.

Рашэнне. У вакне з алфавітам пакідаем толькі кірылічныя літары:

АБВГДЕЖЗИЙКЛМНОПРСТУФХЦЧШЩЫЬЭЮЯ
АбвгдежзийклмнопрстуфхцчшщыьэюяЎўЁёІі.



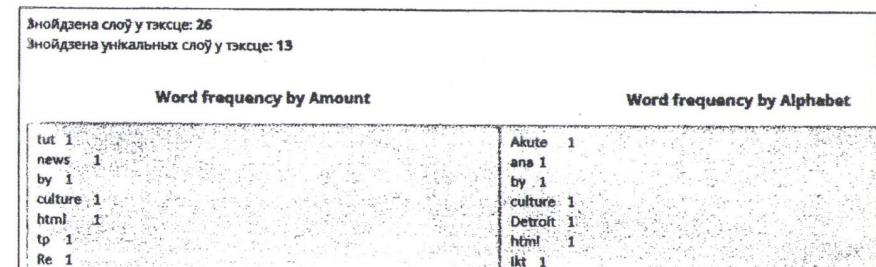
Мал. 3. Рацэнне задачы 2

Сэрвіс [www.corpus.by](#) ўсе слова, якія складаюцца цалкам з кірылічных літар (што, Irap). Але ў выпадку камбінацыі літар і лічбаў апошняя не ўлічваюцца (напрыклад, у слове 45-годдзя знойдзена сімвольная паслядоўнасць годдзя). Таму патрэбна уважліва інтэрпрэтаваць рэзультаты. Аўтары распрацоўкі дапускаюць, што даследчыку спатрэбіца праглядаць вынікі работы сэрвісу і некаторыя выразы, якія не адпавядаюць постаўленай задачы будуть выкрайлены.

Задача 3. Знайсці ўсе слова тэксту, запісаныя лацінскімі літарамі, падлічыць іх частату ўжывання.

Рашэнне. У вакне з алфавітам пакідаем толькі лацінскія літары: ABCDEFGHIJKLMNOPQRSTUVWXYZ abcdefghijklmnopqrstuvwxyz.

Актыўізуем сэрвіс. У выніку маєм (мал. 4):



Мал. 4. Рацэнне задачы 3

У выніку знойдзены ўсе слова, што цалкам складаюцца з лацінскіх літар, але ў выпадку камбінаванага напісання адзінка «разрываеца» на часткі (напрыклад, у назве Re1ikt знойдзены асобныя сімвольныя паслядоўнасці Re i kt). Таму такія слова трэба апрапоўваць асобна.

У бліжэйшай будучыні аўтарамі сэрвісу ставіцца задача дадаць розных мностваў сімвалоў па замоўчванні для шэрагу моў, павялічыць аб'ём магчымай апрапоўкі тэксту. Таксама абмяркоўваюцца пытанні па дапрапоўцы сэрвісу www.corpus.by/wordFrequency з мэтай выкарыстання яго для падліку частотнасці слоў канкрэтнай мовы. Напрыклад, карыстальнік зможа аналізаваць тэксты на грэчаскай мове, змясцішы яе алфавіт у першое мноства, а спіс дыякрытычных сімвалоў – у другое.

ЛІТАРАТУРА

- Frequency of Words [Electronic resource] // www.corpus.by web-site. – 2014. – Mode of access: <http://www.corpus.by/wordFrequency/>. – Date of access: 28.08.2014.
- Go Global Developer Center [Electronic resource] // Windows 1251. – 2014. – Mode of access : <http://msdn.microsoft.com/en-us/goglobal/cc305144.aspx>. – Date of access : 28.08.2014.
- Біблія [Электронны рэсурс] // Беларуская Палічка. – 2014. – Рэжым звароту : <http://knihy.com/none/Biblija.html>. – Дата звароту : 28.08.201