



КАРПОВСКИЕ НАУЧНЫЕ ЧТЕНИЯ

Выпуск 8

Часть I

БЕЛОРУССКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ
ФИЛОЛОГИЧЕСКИЙ ФАКУЛЬТЕТ



КАРПОВСКИЕ НАУЧНЫЕ ЧТЕНИЯ

Сборник научных статей

Основан в 2007 году

Выпуск 8

В двух частях

Часть 1

Минск
«ИВЦ Минфина»
2014

УДК 80/81 (082)

В сборнике представлены материалы Восьмых Карповских научных чтений, посвященных памяти профессора В.А. Карпова — ученого, философа, поэта, оставившего заметный след в отечественной культуре.

Адресован филологам, философам, системологам, специалистам по прикладной и компьютерной лингвистике, а также студентам, магистрантам, аспирантам, интересующимся рассматриваемыми проблемами.

Рекомендовано
Ученым советом филологического факультета
Белорусского государственного университета
(протокол № 9 от 19 июня 2014 г.)

Редакционная коллегия:
кандидат филологических наук, доцент А.И. Головня (отв. ред.);
кандидат филологических наук, доцент Н.С. Касюк (зам. ред.);
кандидат технических наук, доцент О.Е. Елисеева

Рецензенты:
кандидат филологических наук, доцент А.В. Лаврененко
кандидат педагогических наук, доцент Т.В. Игнатович

ISBN 978-985-7060-81-8 (ч. 1)
ISBN 978-985-7060-80-1

© БГУ, 2014
© Оформление.
УП «ИВЦ Минфина», 2014

Секция 3: ОБЩАЯ ТЕОРИЯ СИСТЕМ КАК МЕТОДОЛОГИЯ НАУКИ

Барбук С.Г. (Минск, БГЭУ) Языковые универсалии.....	221
Ван Цзин (Минск, БГУ) Возможность изучения коннотации в корпусном исследовании.....	225
Ван Цин (Минск, БГУ) Особенности имени существительного в грамматике русского и китайского языков.....	228
Гецэвіч Ю.С., Окрут Т.І., Міхайлава Я.А. (Мінск, НАН РБ, БДУ) Распрацоўка лінгвістычных рэсурсаў для алгарытмаў ідэнтыфікацыі рэплік дыялогаў у электронных тэкстах мастацкай тэматыкі на беларускай і рускай мовах.....	231
Гецэвіч Ю.С., Скопінава А.М. (Мінск, АПП НАН Беларусі) Лінгвістычныя рэсурсы для пераўтварэння колькасных выказаў з адзінкамі вымярэння тыпу «лічба-сімвал» у словазлучэнні для беларускай і рускай моў.....	236
Гецэвіч Ю.С. (Мінск, АПП НАН Беларусі), Барадзіна Ю.С. (Мінск, БДУ) Класіфікацыя фразеў дыялогаў па эматыўных прыкметах на матэрыяле рускіх і беларускіх мастацкіх твораў.....	240
Гецэвіч Ю.С., Лысы С.І. (Мінск, АПП НАН Беларусі) Рашэнне прыкладных лінгвістычных задач пры дапамозе сэрвісаў рэсурсу www.corpus.by	243
Глинка Е.В. (Минск, МГЛУ) Проявление симметрии и асимметрии в условиях русско-белорусского двуязычия	247
Головня А.И. (Минск, БГУ) Системная симметрично-асимметричная номинация в аббревиации.....	251
Ивашенко В.П. (Минск, БГУИР) Пространственно-временные интервальные бинарные отношения на множествах событий и их языковые средства представления.....	255

Секция 4: ПРИКЛАДНАЯ ЛИНГВИСТИКА В БЕЛАРУСИ: СОСТОЯНИЕ И ПЕРСПЕКТИВЫ РАЗВИТИЯ

Аскерко Д.С. (Минск, МГЛУ) Специфика базы данных системы автоматического определения средств выражения вербальной агрессии в текстах англоязычных СМИ.....	259
Бочкова А.Л. (Минск, МГЛУ) Лингвистическая база данных как основа системы автоматического извлечения мнений участников интернет-коммуникации.....	263
Гусева Н.Ю. (Минск, БГУ) Трудности в преподавании курса «Основы информационных технологий» иностранным студентам.....	266

**ЛІНГВІСТЫЧНЫЯ РЭСУРСЫ ДЛЯ ПЕРАЎТВАРЭННЯ
КОЛЬКАСНЫХ ВЫРАЗАЎ З АДЗІНКАМІ ВЫМЯРЭННЯ
ТЫПУ «ЛІЧБА–СІМВАЛЬ» У СЛОВАЗЛУЧЭННІ
ДЛЯ БЕЛАРУСКАЙ І РУСКАЙ МОЎ**

Апрацоўка натуральнай мовы з’яўляецца адным з самых актуальных навукова-даследчых накірункаў XXI стагоддзя. Ён прадугледжвае вырашэнне розных камп’ютэрна-лінгвістычных задач [1, с. 333], адной з якіх можна назваць апрацоўку складана ці спецыфічна структураванай інфармацыі ў электронных тэктах: табліц, формул, схем, спасылак, зносак і г.д. У дадзеным дакладзе аўтары канцэнтруюцца на колькасных выразях з адзінкамі вымярэння (КВАВ) — спалучэннях колькасных паказчыкаў (перададзеных на пісьме пасродкам лічбаў) і пазначэнняў мерных адзінак (літарных сімвалаў), напрыклад: $1,72 \text{ г/см}^3$, 19640 км , 55° , $6,387 \text{ сутак}$, $1212 \pm 16 \text{ км}$, $1,9 \times 10^{21} \text{ кг}$ і г.д. Будову КВАВ можна ўявіць фармальна, што адлюстравана ў табліцы 1. У ёй пад X маецца на ўвазе колькасны паказчык (лік), а пад Y — сімвалы літар; знак вертыкальнай рысы размяжоўвае прыклады.

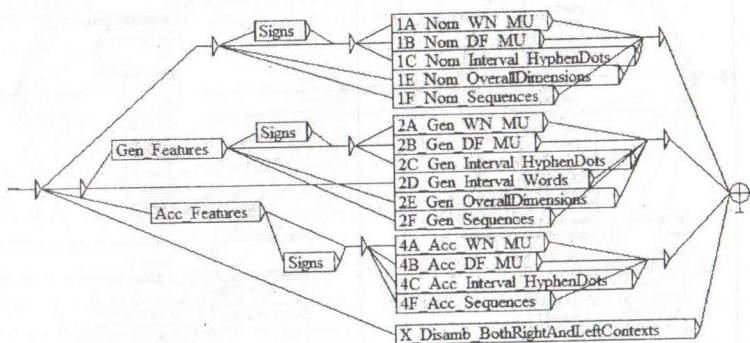
Табліца 1

Мадэлі фармальнага ўяўлення КВАВ

№	Мадэль	Прыклад
1.	$X Y$	12 % 40-50 тыс. м
2.	$X Y/Y$	0,5144444 м/с
3.	$X-[... ..]X Y$	1-1,5 года +13... +19 °С
4.	$X-[... ..]X Y/Y$	$0,1-5,7 \cdot 10^{-2} \text{ м/с}$
5.	$\sim[+\pm><]X Y$	$\pm 0,3^\circ$ $> 6 \text{ Зв}$
6.	$\sim[+\pm\Diamond]X Y/Y$	$\sim 107 \text{ К/с}$ $\sim 9,8 \text{ м/с}^2$
7.	$X, [i] X Y$	2 і 4 метры
8.	$X, X, X Y$	5, 6, 7 шт.
9.	$X Y — X Y$	0,1 Гц — 300 кГц
10.	$X \times X Y$	1136×640 пікселяў
11.	$X \times X \times X Y$	$146,8 \times 75,3 \times 8,9 \text{ мм}$
...

Неабходнасць правільнай апрацоўкі падобных лічба-сімвальных канструкцый актуальная і для навуковай, і для бытавой сферы жыцця, напрыклад, для дадзеных ад штучных спадарожнікаў і касмічных зондаў; медыцынскіх аналізаў (тэмпература, крывяны ціск, пульс, цукар, халестэрын, гемаглабін...); навуковых даследаванняў; кулінарных рэцэптаў; прагнозаў надвор’я; этыкетак на спажывецкіх таварах; апісанняў тавараў у анлайн-крамах; каментарыяў да спартыўных мерапрыемстваў; турыстычных і іншых даведнікаў і г.д.

У працах [2; 3] аўтарамі былі прапанаваныя алгарытмы-рашэнні апрацоўкі КВАВ у выглядзе канчатковых аўтаматаў, якія дазваляюць эмуляваць працу марфалагічных і сінтаксічных граматык. Для іх стварэння быў выкарыстаны наладжвальны лінгвістычны працэсар NooJ. На малюнку 1 прадстаўлена апошняя мадыфікацыя алгарытму, які зараз апрацоўвае КВАВ для беларускай і рускай моў ужо ў трох склонах: назоўным, родным і вінавальным (раней гаворка вялася толькі пра назоўны). Акрамя гэтага, будова дадзенага алгарытму заснаваная ўжо на мадэлях будовы саміх КВАВ, напрыклад, галіна з графам *1E_Nom_OverallDimensions* спрацуе для мадэляў КВАВ пад нумарамі 10 і 11 у табліцы 1.

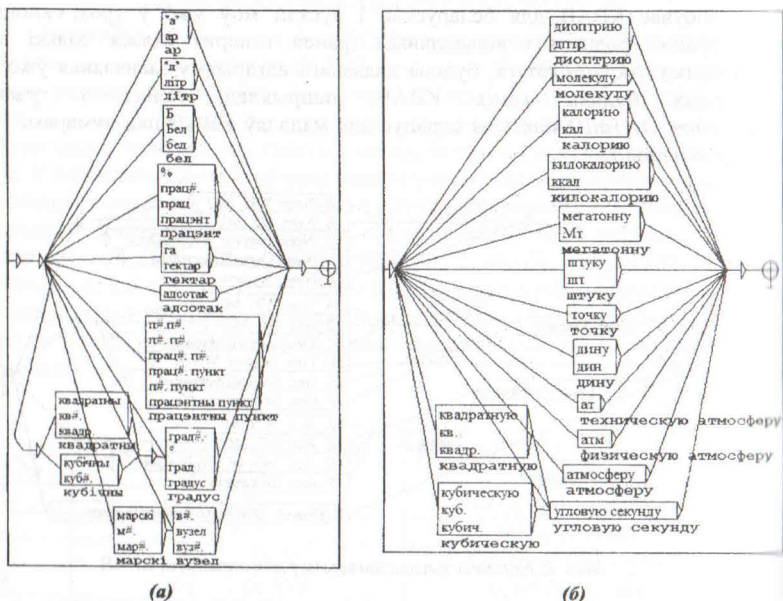


Мал. 1. Агульны выгляд алгарытму для апрацоўкі КВАВ

Апрацоўка КВАВ прадугледжвае іх папярэднюю ідэнтыфікацыю і аналіз. Пад ідэнтыфікацыяй будзем разумець непасрэдны пошук у электронных тэкстах колькасных паказчыкаў (лічбаў), пазначэнняў адзінак (сімвалаў), а таксама спецыфічнага акружэння, якое можа ўплываць на граматычныя катэгорыі ідэнтыфікаванага КВАВ, напрыклад: *5 м* — *5 метраў* (назоўны/вінавальны склон), *каля 5 м* — *каля пяці метраў* (родны склон), *роўны 5 м* — *роўны пяці метрам* (давальны склон) і г.д. Вызначэнне такога ўплыву на КВАВ з боку акружэння і з'яўляецца аналізам. Таксама важна прааналізаваць і элементы КВАВ: прыналежнасць паказчыка колькасці да цэлых ці дзесятковых лікаў, са знакамі ці без знакаў; прыналежнасць ідэнтыфікаванай адзінкі вымярэння да Сістэмы Інтэрнацыянальнай (СИ), распрацаванай Міжнародным бюро мер і вагаў [4]. Нарэшце апрацоўка КВАВ заключаецца ў іх пераўтварэнні ў словы і канчатковым «склеянні» ўсіх элементаў у граматычна правільныя словазлучэнні.

Трэба адзначыць, што нягледзячы на машыннае паходжанне алгарытмаў, іх карэктная праца немагчымая без адпаведных лінгвістычных

рэсурсаў, пад якімі мы маем на ўвазе наборы разгортак колькасных паказчыкаў (лічбаў) і пазначэнняў адзінак вымярэння (літарных сімвалаў) у цэлыя словы. Іх распрацоўка і была задачай дадзенага дакладу (малюнак 2).



Мал. 2. Лінгвістычны рэсурсы ў выглядзе графа для апрацоўкі пазначэнняў дадатковых адзінак вымярэння для беларускай (а) і рускай (б) моў

Падкрэслім, што лінгвістычныя рэсурсы непасрэдна ўбудаваныя ў алгарытм у выглядзе асобных графаў, а не падключаныя, напрыклад, у якасці самастойных слоўнікаў. На дадзены момант апрацоўваецца 120 адзінак вымярэння ў розных відах запісу. Дзеля зручнасці стварэння і паўнаты рэсурсаў патрэбна было перш за ўсё скласці ліст з адзінкамі вымярэння і размежаваць адзінкі ў некалькі груп. За аснову была выкарыстаная класіфікацыя Міжнароднага бюро мер і вагаў: базавыя адзінкі СІ, вытворныя ад адзінак СІ і пазасістэмныя адзінкі. Пазней, па меры таго як алгарытм увесь час тэставалася на беларуска- і рускамоўных тэкставых масівах навукова-тэхнічнай тэматыкі, спіс адзінак значна папоўніўся. Так, усе астатнія адзінкі, не апісаныя СІ, а таксама словы, якія ўмоўна можна лічыць адзінкамі вымярэння (напрыклад, *шт.* ад *штука*), увайшлі ў алгарытм асобным дадатковым класам Extra. На малюнку 2 як раз прадстаўлены выгляд гэтага графу менавіта для беларускай і рускай моў. Вынікі выкарыстання алгарытма дэманструюцца ў табліцы 2.

Фрагменты вынікаў апрацоўкі алгарытмам беларуска- (а) і рускамоўнай (б) навукова-тэхнічнай тэкставай інфармацыі

Выгляд у тэксце	Пасля апрацоўкі алгарытмам
(а)	
займае 36,4 % 1136x640 пікселяў на 1,3 мегапікселя на 1,5-7 адсоткаў 1430 МАг 500-600 кв. метраў	займае трыццаць шэсць цэлых чатыры дзясятыя працэнта адна тысяча сто трыццаць шэсць на шэсцьсот сорок пікселяў на адну цэлую тры дзясятыя мегапікселя на адну цэлую пяць дзясятых дэфіс сем адсоткаў адна тысяча чатырыста трыццаць міліампер-гадзін пяцьсот дэфіс шэсцьсот квадратных метраў
(б)	
0,01 МДж ~ 0,1 МДж болей 40 га 20-210 мм +25° С свыше 2000 м	ноль цэлых адна сотая мегаджоуля около нуля цэлых одной десятой мегаджоуля болей сорока гектаров двадцать дефис двести десять миллиметров плюс двадцать пять градусов Цельсия свыше двух тысяч метров

У заключэнне падкрэслім: значнасць лінгвістычных рэсурсаў для камп'ютэрных рашэнняў задач, звязаных з апрацоўкай тэкставай інфармацыі, немагчыма пераацаніць. Чым лепш распрацаваныя рэсурсы, тым вышэйшае значэнне паўнаты ўсяго алгарытму — найбольш важнага паказчыка якасці алгарытму, акрамя дакладнасці. Паўната вылічваецца як вынік дзялення колькасці КВАВ, якія алгарытм правільна апрацаваў, на рэальную колькасць КВАВ ва ўсім тэкставым мностве, якую падлічыў эксперт. Зараз паўната алгарытму дасягае 75 %. У бліжэйшых планах аўтараў палепшыць гэтае значэнне праз тэставанне алгарытма і наступную дапрацоўку яго лінгвістычных рэсурсаў на дадатковым тэкставым мностве. У доўгатэрміновай перспектыве мае быць укараненне алгарытму ў сістэму сінтэзу маўлення па тэксце дзеля большай сэнсавай дакладнасці і правільнасці, а таксама для правільнай інтанацыйнай і прасадычнай афарбоўкі тэкстаў для агучвання.

ЛІТАРАТУРА

1. Гецэвіч, Ю.С. Метад пабудовы кампанентаў сінтэзу маўлення па тэксце для натуральна-маўленчага інтэрфейса пры дапамозе NooJ / Ю.С. Гецэвіч, А.М. Скопінава, Т.І. Окрут // Открытые семантические технологии проектирования интеллектуальных систем = Open Semantic Technologies for Intelligent Systems (OSTIS-2014): материалы IV Междунар. науч.-техн. конф. (Минск, 20–22 февраля 2014 года) / редкол.: В.В. Голенков (отв. ред.) [и др.]. – Минск: БГУИР, 2014 г. – С. 333–338.
2. Hetsevich, Yu.S. Transforming quantitative expressions with measurement units into orthographical words for text-to-speech synthesis to Belarusian and Russian / Yu.S. Hetsevich, A.M. Skopinava // Вестник МГУЛУ. Сер. 1, Филология. – 2013. – № 3. – С. 133–144.
3. Гецэвіч, Ю.С. Мадэляванне і распрацоўка сістэм пошуку колькасных выразаў з адзінкамі вымярэння ў электронных тэкстах на беларускай і рускай мовах / Ю.С. Гецэвіч,

А.М. Скопінава, А.Ф. Есіс // Развитие информатизации и государственной системы научно-технической информации (РИНТИ-2013): доклады XII Международной конференции (Минск, 20 ноября 2013 г.). – Минск: ОИПИ НАН Беларуси, 2013. – С. 282–287.

4. Апісанне СІ на сайце Міжнароднага бюро мер і вагаў [Электронны рэсурс]. – 2006. – Рэжым доступу: http://www.bipm.org/en/si/si_brochure/general.html. – Дата доступу: 15.04.2014.

Ю.С. Гецэвіч (Мінск, АПН НАН Беларусі), Ю.С. Барадзіна (Мінск, БДУ)

КЛАСІФІКАЦЫЯ ФРАЗАЎ ДЫЯЛОГАЎ ПА ЭМАТЫЎНЫХ ПРЫКМЕТАХ НА МАТЭРЫЯЛЕ РУСКІХ І БЕЛАРУСКІХ МАСТАЦКІХ ТВОРАЎ

Сінтэз маўлення знаходзіць прымяненне ў розных сферах, напрыклад, пры стварэнні аўтаадказчыкаў, натуральна-моўных інтэрфейсаў, агучванні інфармацыйных паведамленняў у транспарце, на вакзале, у аэрапорце і г.д. Акрамя таго, тэхналогіі сінтэзу маўлення могуць выкарыстацца для машыннага стварэння аўдыёкніг. Звычайна гэта вымагае сур'езнай працы дыктараў і актораў, але з дапамогай існуючых напрацовак у сферы сінтэзу маўлення працэс можа быць аўтаматызаваны.

У лабараторыі распазнавання і сінтэза маўлення Аб'яднанага інстытута праблем інфарматыкі НАН Беларусі быў распрацаваны сінтэзатар, які тэхнічна ўжо можа “чытаць кнігі” [1, 269]. Тым не менш, працэс якаснага сінтэзу маўлення яшчэ не скончаны, і застаецца некалькі істотных праблемаў, якія дагэтуль не былі вырашаны, і адна з іх — гэта праблема інтанацыі.

У п'есах аўтары спрашчаюць працу актораў з дапамогай рэмарак, якія падказваюць неабходную інтанацыю. У дыялогах праявілі тэкстаў прысутнічаюць словы аўтара, здольныя выконваць гэтую ж функцыю. Так, прыведзены ніжэй тэкст, калі будзе агучаны акторм у аўдыёкнізе, ніколі не будзе прачытаны манатонна: эмоцыі моўцы бачны праз знакі прыпынку, сутнасць самой фразы, і, не ў апошнюю чаргу, праз «падказкі» ў словах аўтара.

— Ад нараджэння вольныя! — люта роў ён. — Вось вам вашы вольнасьці! Усіх іх выразаць!

Так, для якаснага стварэння аўдыёкніг з дапамогай сінтэзатара маўлення трэба ўлічваць эмоцыі, закладзеныя ў рэпліках герояў і словах аўтара.

Такім чынам, наступны артыкул знаёміць з першаснымі вынікамі працы, мэта якой была ў тым, каб знайсці ў фразях дыялогаў ідэнтыфікатары эматыўнай палярнасці (пазітыўны, негатыўны, нейтральны) і прапанаваць магчымыя сродкі іх фармалізацыі, а таксама працэс іраваць іх на невялікім працоўным корпусе.

Даследаванне вядзецца адначасова для беларускай і рускай моў. Для беларускай мовы ў якасці матэрыяла быў выбраны твор Алеся Рукала «На