

XII Международная конференция

**РАЗВИТИЕ ИНФОРМАТИЗАЦИИ
И ГОСУДАРСТВЕННОЙ СИСТЕМЫ
НАУЧНО-ТЕХНИЧЕСКОЙ
ИНФОРМАЦИИ**

РИНТИ-2013



20 ноября 2013 года, Минск

ДОКЛАДЫ

УДК 002; 004

• **Развитие информатизации и государственной системы научно-технической информации (РИНТИ-2013)** : доклады XII Международной конференции, Минск, 20 ноября 2013 г. – Минск : ОИПИ НАН Беларуси, 2013. – 412 с. – ISBN 978-985-6744-78-8.

Представлены доклады XII Международной конференции «Развитие информатизации и государственной системы научно-технической информации» (РИНТИ-2013), Минск, 20 ноября 2013 г., в которых рассматриваются вопросы научно-методического, информационного и технологического обеспечения развития информатизации, создания и развития автоматизированных систем научно-технической информации, корпоративных библиотечно-информационных систем и технологий, а также направления использования информационно-коммуникационных технологий и результаты научных исследований в данных областях.

Материалы конференции будут полезны специалистам в области информационно-коммуникационных технологий, занимающимся разработкой и внедрением автоматизированных информационных систем управления, развитием информационной инфраструктуры Беларуси, реализацией проектов государственных программ в сфере информатизации и систем научно-технической информации.

Одобрены программным комитетом и печатаются по решению редакционной коллегии Объединенного института проблем информатики Национальной академии наук Беларуси в виде, представленном авторами.

Научные редакторы:

доктор физико-математических наук, профессор А.В. Тузиков;
кандидат технических наук, доцент Р.Б. Григянец;
кандидат технических наук В.Н. Венгеров

ISBN 978-985-6744-78-8

© ГНУ «Объединенный институт
проблем информатики Национальной
академии наук Беларуси», 2013

Григянец Р.Б., Венгеров В.Н., Мисякова Г.Т., Лаужель Г.О. Системы и технологии автоматизации библиотечной и информационной деятельности	238
Горбачев Н.Н. Технологии параллельной визуализации информации	247
Горбачев Н.Н. Сценарные методы в работе ситуационно-аналитических центров	253
Липницкий С.Ф. Анализ тонально-окрашенных текстовых сообщений в Интернете	259
Токарь О.В. Семантический профиль областей экрана для разработки электронных изданий	266
Гецэвіч Ю.С., Округт Т.І., Лабанаў Б.М. Аўтаматызацыя шматгаласавога стварэння аўдыёкніг на беларускай мове з дапамогай сінтэзатораў маўлення па тэксле	269
Гецэвіч Ю.С., Пакладок Д.А., Брэк Д.В. Сінтэзатар маўлення па тэксле для шырокага кола карыстальнікаў Інтэрнэта	277
Гецэвіч Ю.С., Скопінава А.М., Есіс А.Ф. Мадэляванне і распрацоўка сістэм пошуку колькасных выразаў з адзінкамі вымярэння ў электронных тэкстах на беларускай і рускай мовах	282
Гецэвіч Ю.С., Кожух В.С., Пракаповіч Р.А., Сычоў У.А., Герасюта С.Л. Распрацоўка адкрытага праграмнага забеспячэння модульнай платформы для канструявання і кіравання навучальна-даследчымі робататэхнічнымі апаратамі	288
Лаужель Г.О., Молчан Ж.М., Степанцова Е.В. Автоматизированная информационная система в области экологии, окружающей среды и природопользования	294
Ильина З.М., Кондратенко С.А., Гридюшко Д.Н., Мисякова Г.Т. Автоматизированная система информационного обеспечения инновационной деятельности на национальном рынке сельскохозяйственного сырья и продовольствия	300
Сафонов Р.Ф., Лаужель Г.О. Автоматизированная система учета и формирования ведомственной отчетности о причинах временной нетрудоспособности в учреждениях НАН Беларуси	306
Петухов А.В. Распределение ролей пользователей типовой системы профессионального образования в области разработки и внедрения интегрированных систем проектирования и производства	312
Ганченко В.В., Дудкин А.А., Петровский А.И. Информационные технологии в задаче мониторинга состояния сельскохозяйственной растительности	315

МАДЭЛЯВАННЕ І РАСПРАЦОЎКА СІСТЭМ ПОШУКУ КОЛЬКАСНЫХ ВЫРАЗАЎ З АДЗІНКАМІ ВЫМЯРЭННЯ Ў ЭЛЕКТРОННЫХ ТЭКСТАХ НА БЕЛАРУСКАЙ І РУСКАЙ МОВАХ

Ю.С. Гецэвіч¹, А.М. Скопінава¹, А.Ф. Есіс²

¹Аб'яднаны інстытут праблем інфарматыкі НАН Беларусі, Мінск;

²Беларускі дзяржаўны ўніверсітэт інфарматыкі і радыёэлектронікі, Мінск

Разгледжсаны працэс пабудовы алгарытмічных мадэляў і эксперыментальна-праграмнага комплексу для вырашэння задачы пошуку і класіфікацыі колькасных выразаў з адзінкамі вымірэння на матэрыяле беларуска- і рускамоўных электронных тэкстаў.

Уводзіны

Электронныя тэксты навукова-тэхнічнай тэматыкі часта ўтрымліваюць колькасныя выразы з адзінкамі вымірэння (КВАВ), напрыклад 27 %, 110 м/с, 15,25 кг і інш. Яны павінны быць дакладна пабудаваны і ўжыты, каб карыстальнік тэксту мог правільна і адназначна зразумець колькасны бок пэўнага якаснага паказчыка.

Падчас апрацоўкі колькасных выразаў з адзінкамі вымірэння ўзнікаюць разнастайныя задачы і праблемы ў наступных сферах [1–3]: *выдавецкія ўстановы* (проблема аўтаматызаванай лакалізацыі канкрэтнага спісу выразаў з адзінкамі вымірэння і задача хуткай праверкі правільнасці ўжывання разгорнутых формаў называў адзінак вымірэння ў тэкстах); *сістэмы сінтэзу маўлення па тэксце* (проблема правільнай генерацыі арфаграфічнага тэксту па ўваходным тэксце); *сістэмы пошуку і апрацоўкі інфармацыі* для бібліятэк, Інтэрнэту, патэнтавых базаў дадзеных (задача фармавання дакладных пошукавых запытаў; задачы аўтаматычнага рэферавання і анафавання інфармацыі).

У дакладзе разгледжаны абагульнены падыход да решэння камп'ютарна-лінгвістычных задач на электронных тэкстах на прыкладзе задачы пошуку колькасных выразаў з адзінкамі вымірэння ў навукова-тэхнічных тэкстах.

1. Агульны падыход да решэння камп'ютарна-лінгвістычных задач па электронных тэкстах

Пад камп'ютарна-лінгвістычнай задачай па электронным тэксле будзем разумець такую задачу (проблему), якая, па-першае, ставіцца адносна электроннага тэксту, па-другое, тычыцца пытанняў канкрэтнага пошуку, класіфікацыі ці перапрацоўкі паслядоўнасці электронных сімвалоў мэтавага электроннага тэксту, па-трэцяе, яе канчатковым решэннем павінна быць камп'ютарная програма, праца якой можа быць праверана карыстальнікам на неабмежаванай колькасці іншых электронных тэкстаў. Такое азначэнне з'яўляецца абагульненым падыходам да шэрагу задач, якія ставіліся і вырашаліся ў працах [1–3, 5]. Прагледзім працэс ад пастаноўкі канкрэтнай задачы да яе вырашэння з улікам камп'ютарных сродкаў і ўмоваў. На мал. 1 бачна, што дадзены працэс прадугледжвае ажыццяўленне шасці этапаў ад вызначэння задачы да стварэння прадукту (яе вырашэння) для карыстальніка.

Такім чынам, спачатку перад даследчыкам ставіцца праблема (1). Пасля эксперыментальнымі шляхам эксперт заходзіць решэнні праблемы адносна невялікага

фрагменту тэксту (2), для гэтых распрацоўваеца дзейсны мадэль-алгарытм з дапамогай канчатковых аўтаматаў NooJ (3) [4]. Далей распрацаваныя канчатковыя аўтаматы тэстуюцца на большай колькасці тэкстаў (4).



Мал. 1. Працэс вырашэння камп'ютарна-лінгвістычнай задачы

Калі вынікі тэставання паводле ацэнак дакладнасці і паўнаты не задавальняюць, то мадэль-алгарытм вяртаецца для дапрацоўкі на стадью (2), у супрацьлеглым выпадку ён перадаецца для стварэння на яго базе эксперыментальна-праграмнага комплексу (5). Прычым тэставаныя дадзеныя, якія былі распрацаваны на этапе (4), перадаюцца на этапы (5) і (6) для дакладнай распрацоўкі і тэставання праграмнага прадукта. Пасля распрацоўкі праграмы (5) адбываюцца яе тэставанне і ацэнка дакладнасці (6). Пры становчай ацэнцы (ў межах дапушчальнай памылкі) праграма трапляе ў рукі да карыстальніка (7) і набывае статус канчатковага прадукту, а пры адмоўнай – адбываеца вяртанне на стадью (5).

2. Пошук колькасных выразаў з адзінкамі вымярэння як камп'ютарна-лінгвістычная задача

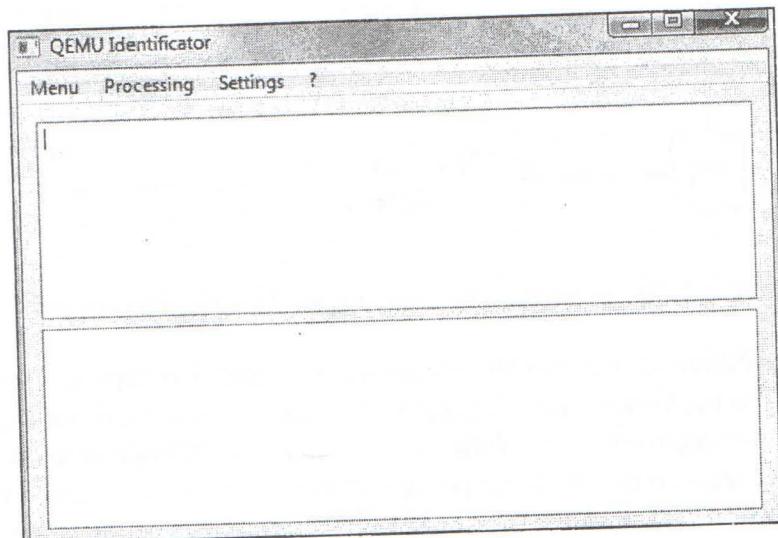
Спачатку паводле схемы на мал. 1 была вызначана канкрэтная задача (1): знайсці і класіфікаціаць у электронных тэкстах КВАВ. На этапе (2) экспертам праводзіліся назіранні і аналіз колькасных выразаў з адзінкамі вымярэння адносна іх будовы і выкарыстання на мэтавым матэрыяле электронных тэкстаў на беларускай і рускай мовах. На дадзены момант паводле этапу (3) аўтарамі артыкулу змадэляваны трох ўзаемадапаўняючыя алгарытмы для пошуку КВАВ у вялікіх корпусах тэкстаў, якія дазваляюць:

- знаходзіць КВАВ і класіфікаць іх па трох тыпах паводле Міжнароднай сістэмы адзінак CI (CI, вытворныя ад CI, не CI) [2];
- знаходзіць КВАВ з метралагічнымі прыстаўкамі (кратнымі ці дольнымі, скарочанымі ці ў поўнай форме) і класіфікаць іх паводле словаўтваральныхых асаблівасцяў [3];
- пераўтвараць КВАВ у арфаграфічныя слоўкі [1, 3].

Этап тэставання (4) паказаў, што, напрыклад, першы алгарытмічны комплекс дае пошукавыя вынікі з дакладнасцю ў 72 %. Гэты адносна высокі паказчык даў падставу для распрацоўкі эксперыментальнага праграмнага комплексу (5), які б знаходзіў КВАВ і класіфікаць іх па трох тыпах адносна Міжнароднай сістэмы адзінак CI (CI, вытворныя ад CI, не CI).

3. Рэалізацыя эксперыментальнай праграмы пошуку і ідэнтыфікацыі колькасных выразаў з адзінкамі вымярэння паводле Міжнароднай сістэмы адзінак

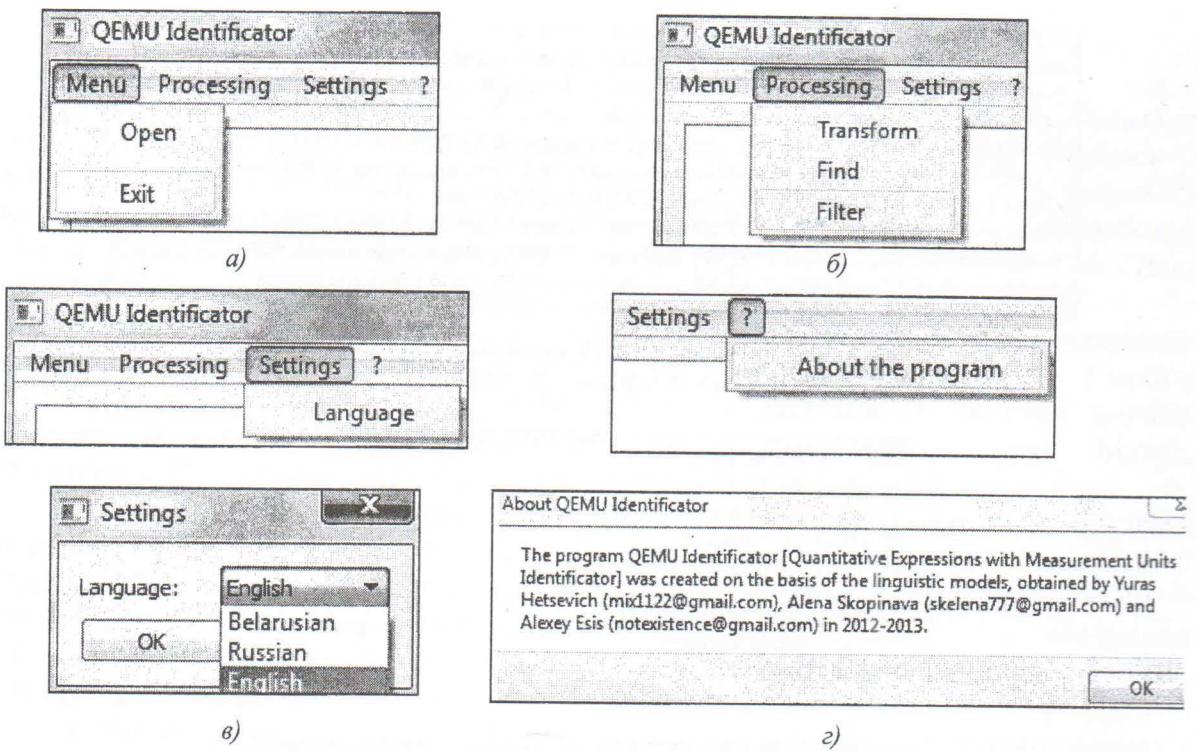
На мал. 2 дэманструеца запуск праграмы з называй *QEMU Identifier (Quantitative Expressions with Measurement Units Identifier – Ідэнтыфікатор колькасных выразаў з адзінкамі вымярэння)*.



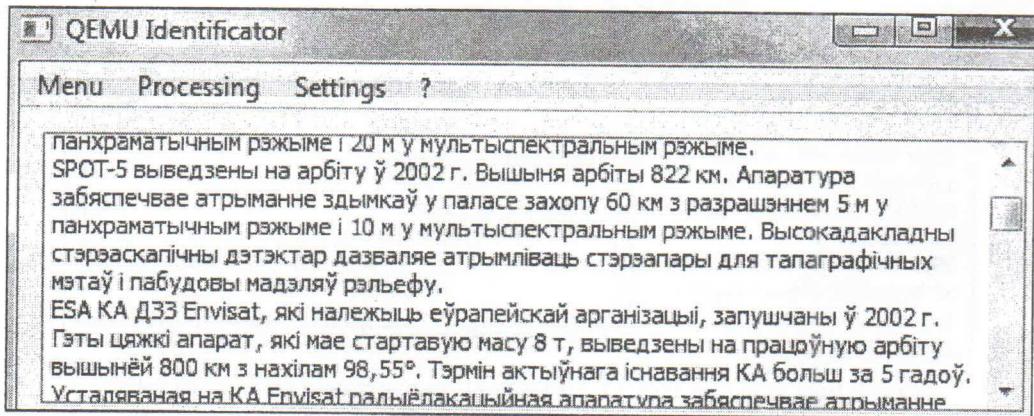
Мал. 2. Запуск эксперыментальна-праграмнага комплексу *QEMU Identifier*

Панэль кіравання складаецца з наступных укладак: *Menu* (Меню), *Processing* (Апрацоўка), *Settings* (Налады) і (даведка) (мал. 3). Праз меню карыстальнік можа адкрыць файл (*Open*) і выйсці з праграмы (*Exit*) (мал. 3, а). Укладка апрацоўкі прадугледжвае аналіз тэксту праз каманды пераўтварыць (*Transform*), знайсці (*Find*) і фільтраваць (*Filter*) (мал. 3, б). Праз укладку налад даецца магчымасць выбараць мову інтэрфейсу (*Language*): беларускую (*Belarusian*), рускую (*Russian*) або англійскую (*English*) (мал. 3, в). Нарэшце ў даведцы падаюцца агульная інфармацыя аб праграме і контактныя дадзеныя распрацоўшчыкаў (*About the program*) (мал. 3, г).

Разгледзім працу *QEMU Identifier* на прыкладзе апрацоўкі фрагменту тэксту з навукова-тэхнічнага беларускамоўнага тэкставага корпусу. Тэксты для аналізу можна ўводзіць адвольна ў верхнє поле (мал. 2), а можна адкрываць гатовыя тэкставыя фрагменты ў фармаце txt праз адпаведную каманду (мал. 3, а). У нашым выпадку возьмем гатовы (мал. 4).



Мал. 3. Знаменства з інтэрфейсам эксперыментальна-праграмнага комплексу *QEMU Identifier*



Мал. 4. Фрагмент навукова-тэхнічнага тэксту для аналізу праз *QEMU Identifier*

Прадэманструем, што адбываеца пры выкарыстанні каманд:

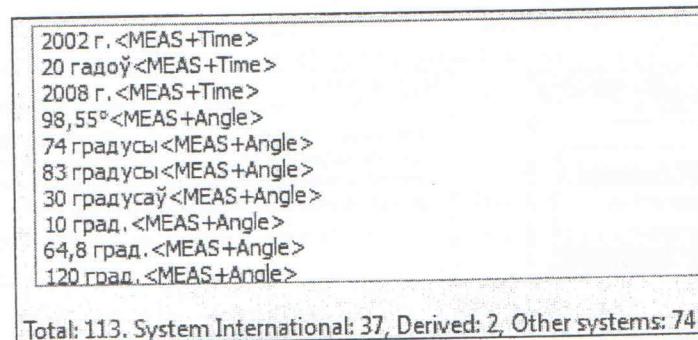
Transform – знайдзеным КВАВ надаюца пэўныя маркеры, вынік адлюстроўваеца ў ніжнім полі інтэрфейсу; да прыкладу: у тэксле сустрэлася *15 m*, праграма яго пераўтварыла ў *{15 m}<MEAS+Length|Distance+SI>* – гэта азначае, што дадзены літарна-сімвальны набор з'яўляеца колькасным выразам з адзінкай вымярэння даўжыні або адлегласці і гэта адзінка належыць сістэме СІ (мал. 5);

Find – выводзіца спіс усіх знайдзеных КВАВ і ніжэй падаюца статыстычныя звесткі: агульная колькасць КВАВ, іх разнастайнасць; да прыкладу: у тэставым тэксле знайшлося 113 КВАВ, з якіх 37 выразаў змяшчаюць сістэмныя мерныя адзінкі, 2 – вытворныя і 74 – пазасістэмныя (мал. 6);

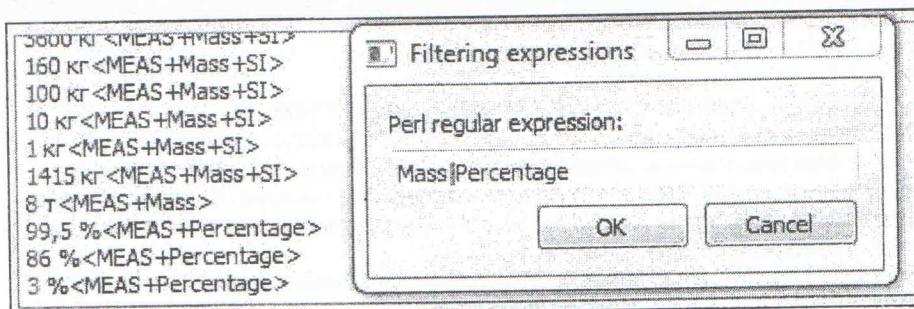
Filter – даеца магчымасць здзяйсняць разнастайныя пошукавыя запыты праз фармальную мову рэгулярных выразаў; напрыклад, праз выраз *Mass|Percentage* можна атрымаць адразу спіс КВАВ і масы, і працэнта (мал. 7).

Амерыканская праграма Landsat з'яўляецца адной з найбольш прадуктыўных на сусветным рынку дадзеных ДЗЗ, з {1972 г.} <MEAS+Time> у межах праграмы было паспяхова запушчана шэсць спадарожнікаў. Дзеісны на сёння спадарожнік Landsat-7 знаходзіцца на арбіце з {1999 г.} <MEAS+Time>. На спадарожніку ўсталяваны радыёметр ETM+, які мае панхраматычны канал высокага разрашэння ({15 м} <MEAS+Length|Distance+SI>). Іншы амерыканскі КА Terra быў запушчаны ў снежні {1999 г.} <MEAS+Time> у межах праграмы EOS (Earth Observing System). Абсталяваны апаратным комплексам дыстанцыйнага зандавання ASTER, які выкарыстоўваецца для дэталёвага

Мал. 5. Маркіраванне КВАВ у навукова-тэхнічным тэксле праз каманду *Transform* у *QEMU Identifier*



Мал. 6. Вынікі пошуку КВАВ у навукова-тэхнічным тэксле праз каманду *Find* у *QEMU Identifier*



Мал. 7. Вынікі фільтрацыі КВАВ масы ці працэнта праз рэгулярны выраз у навукова-тэхнічным тэксле праз каманду *Filter* у *QEMU Identifier*

Заключэнне

Разгледжаны абагульнены падыход да решэння камп'ютарна-лінгвістычных задач па электронных тэкстах на прыкладзе задачы пошуку колькасных выразаў з адзінкамі вымярэння ў навукова-тэхнічных тэкстах.

Распрацаваны праграмны эксперыментальны комплекс дазваляе зручна шукаць колькасныя выразы з адзінкамі вымярэння сістэмы СІ і вытворных ад СІ у электронных тэкстах на беларускай ці рускай мовах. Зараз гэтая реалізацыя тэстуеца распрацоўшчыкамі і карыстальнікамі. На базе вынікаў, якія будуть атрыманы, плануецца ў далейшым палепшыць гэты праграмны прадукт, а таксама праграмна реалізуваць астатнія алгарытмічныя мадэлі для пошуку колькасных выразаў з адзінкамі вымярэння.

Спіс літаратуры

1. Skopinava, A.M. Processing of quantitative expressions with units of measurement in scientific texts as applied to Belarusian and Russian text-to-speech synthesis / A.M. Skopinava, Yu.S. Hetsevich, B.M. Lobanov // Комп'ютерная лингвистика и интеллектуальные технологии : материалы Междунар. конф. «Диалог», Московская обл., г. Бекасово, 29 мая – 2 июня 2013 г. – Вып. 12 (19). – В 2 т. – Т.1. – М. : Изд-во РГГУ, 2013. – С. 634–651.
2. Гецэвіч, Ю.С. Ідэнтыфікацыя выразаў з адзінкамі вымярэння ў навукова-тэхнічных і прававых тэкстах на беларускай і рускай мовах / Ю.С. Гецэвіч, А.М. Скопінава // Развитие информатизации и государственной системы научно-технической информации (РИНТИ-2012) : доклады XI Междунар. конф., Минск, 15 ноября 2012 г. – Минск : ОІПІ НАН Беларусі, 2012. – С. 260–265.
3. Гецэвіч, Ю.С. Кампаненты ідэнтыфікацыі колькасных выразаў з адзінкамі вымярэння ў тэкстах на беларускай і рускай мовах / Ю.С. Гецэвіч, А.М. Скопінава // Открытые семантические технологии проектирования интеллектуальных систем = Open Semantic Technologies for Intelligent Systems (OSTIS-2013) : материалы III Междунар. науч.-техн. конф., Минск, 21–23 февраля 2013 г. – Минск : БГУИР, 2013 г. – С. 319–328.
4. Лінгвістычны працэсар NooJ [Электронны рэсурс]. – 2002. – Рэжым доступу : <http://www.nooj4nlp.net/pages/nooj.html>. – Дата доступу : 01.07.2012.
5. Гецэвіч, Ю.С. Аўтаматызаваная апрацоўка сімвальных выразаў у тэкстах для сістэмы сінтэзу беларускага маўлення / Ю.С. Гецэвіч // Информатика. – 2011. – № 4. – С. 82–93.