



OSTIS-2015

(Open Semantic Technologies for Intelligent Systems)

УДК 004.934.5

АДКРЫТЫЯ КАМПАНЕНТЫ WWW.CORPUS.BY ДЛЯ НАТУРАЛЬНА-МАЎЛЕНЧАГА ІНТЭРФЕЙСУ

Гецэвіч Ю.С.*, Лабанаў Б.М.*, Лысы С.І.*, Гюнтар А.В.*, Дзенісюк Д.А.*, Захар'еў В.А.**

** Аб'яднаны інстытут праблем інфарматыкі Нацыянальнай акадэміі навук Беларусі,
г. Мінск, Беларусь*

yury.hetsevich@gmail.com

lobanov@newman.bas-net.by

stanislau.lysy@gmail.com

lena205593@gmail.com

d.denissyuk@gmail.com

*** Беларускі дзяржаўны ўніверсітэт інфарматыкі і радыёэлектронікі, г. Мінск, Беларусь*

delfvad@gmail.com

У дадзеным артыкуле апісаны агульныя падыходы да праектавання, рэалізацыі і выкарыстання адкрытых кампанентаў інтэрнэт-рэсурсу www.corpus.by для натуральна-маўленчага інтэрфейсу з мэтай паляпшэння і ўдасканалення працы сістэмы сінтэзу маўлення па тэксце з адначасовым прадстаўленнем магчымасці інтэграцыі з праектам OSTIS. Прыводзяцца апісанні мэтаў і задачай, спосабаў працы і даступнай функцыянальнасці некаторых рэалізаваных сэрвісаў.

Ключавыя словы: сінтэз маўлення па тэксце, натуральна-маўленчы інтэрфейс, кампанент, інтэрнэт-сэрвіс.

Уводзіны

Дзеля вырашэння задачы ўсталявання ўзаемадзеяння паміж чалавекам і машынай пры дапамозе маўлення навукоўцы ва ўсім свеце працуюць над стварэннем сістэм сінтэзу і распазнавання маўлення. У апошні час сістэмы сінтэзу маўлення па тэксце (ССМТ) ужо дасягнулі пэўнай дасканаласці і шырока выкарыстоўваюцца. Аднак нельга сказаць, што задача сінтэзу маўлення з'яўляецца вырашанай, асабліва ў дачыненні да беларускай мовы. У сувязі з гэтым аўтарамі было прынята рашэнне стварыць інтэрнэт-рэсурс, які б аб'ядноўваў шэраг інтэрнэт-сэрвісаў, кожны з якіх у сваю чаргу вырашаў бы пэўныя свае ўласныя задачы (напрыклад, прыкладныя лінгвістычныя задачы), але ў той жа час дапамагаў распрацоўшчыкам ва ўдасканаленні працы сінтэзатара маўлення па тэксце. Такі падыход дае магчымасць праз ужыванне сэрвісаў вялікай колькасцю карыстальнікаў правесці тэставанне і збор неабходнай інфармацыі, пры дапамозе якіх можа ўдасканальвацца якасць працы сінтэзатара маўлення па тэксце [Лобанов, 2008], [Гецэвіч, 2012].

У дакладзе прадстаўлены агульныя падыходы да праектавання, рэалізацыі і выкарыстання інтэрнэт-сэрвісаў у мэтах паляпшэння працы ССМТ з прадстаўленнем магчымасці інтэграцыі з іншымі праектамі (у прыватнасці з праектам OSTIS) праз інтэрнэт-сайт www.corpus.by [www.corpus.by, 2012]. Прыводзяцца апісанні мэтаў і задачай, спосабаў працы і даступнай функцыянальнасці некаторых рэалізаваных сэрвісаў.

1. Агульныя падыходы да праектавання інтэрнэт-сэрвісаў

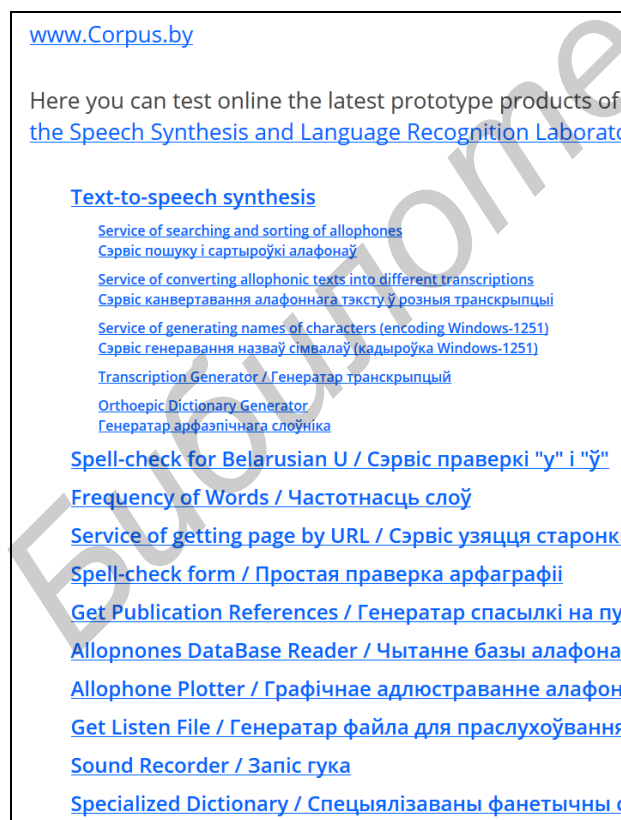
Распрацоўка ССМТ з'яўляецца складанай задачай, якая ахоплівае вялікую колькасць падзадач, а таксама сутыкаецца з шэрагам пабочных прыкладных лінгвістычных задач, якія патрабуюць вырашэння. Задачы, якія паўстаюць пры стварэнні ССМТ, можна ўмоўна падзяліць на тры групы:

- перадапрацоўку ўваходных тэкстаў;
- непасрэдна сінтэз маўлення;
- апрацоўку і выкарыстанне выніковых даных.

Інтэрнэт-рэсурс www.corpus.by мае сваёй мэтай рашэнне як непасрэднай задачы сінтэзавання

маўлення, так і шэрагу дадатковых задач перадапрацоўкі тэксту і апрацоўкі вынікаў працы сінтэзатара маўлення. Дзеля дасягнення дадзенай мэты было прынята рашэнне распрацаваць для гэтага рэсурсу сістэму інтэрнэт-сэрвісаў, кожны з якіх вырашаў бы пэўную ўласную задачу, а ў спалучэнні з іншымі сэрвісамі – неабходныя падзадачы ССМТ. На сённяшні дзень распрацаваны альбо знаходзяцца ў праежжавай стадыі распрацоўкі і даступны бясплатна ў Інтэрнэце наступныя сэрвісы:

- сэрвіс пошуку і сартыроўкі алафонаў;
- сэрвіс канвертавання алафоннага тэксту ў розныя транскрыпцыі;
- сэрвіс генерацыі назваў сімвалаў;
- сэрвіс генерацыі транскрыпцый;
- сэрвіс генерацыі арфаэпічнага слоўніка;
- сэрвіс праверкі “y” і “ў”;
- сэрвіс падліку частотнасці слоў;
- сэрвіс генерацыі спасылкі на публікацыю;
- сэрвіс графічнага адлюстравання алафонаў і алафонных фраз;
- сэрвіс генерацыі файла для праслухоўвання;
- сэрвіс запісу гуку;
- сэрвіс спецыялізаванага фанетычнага слоўніка; і іншыя.



Малюнак 1 – Знешні інтэрфейс інтэрнэт-рэсурсу www.Corpus.by

Знешні інтэрфейс пачатковай старонкі інтэрнэт-рэсурсу www.Corpus.by, на якой знаходзяцца

спасылкі на ўсе пералічаныя вышэй сэрвісы, прадстаўлены на малюнку 1.

2. Фарміраванне патрабаванняў да праектуемых інтэрнэт-сэрвісаў

Мэтай распрацоўкі інтэрнэт-сэрвісаў з’яўляецца павышэнне якасці працы сінтэзатара маўлення па тэксце для беларускай мовы. Аўтарамі артыкула было заўважана, што гэтага можна дасягнуць некалькімі спосабамі: па-першае, праз вынясенне пэўных элементаў існуючай ССМТ длявольнага карыстання ў Інтэрнэце у выглядзе інтэрнэт-сэрвісаў, якія б мелі пэўную мэтавую аўдыторыю; па-другое, праз распрацоўку пэўных інтэрнэт-сэрвісаў, якія б маглі быць у далейшым інтэграваны ў ССМТ, а пакуль удасканалення і адтэставання ў Інтэрнэце. Такім чынам можна атрымаць наступнае:

- непасрэдня вынікі працы кожнага сэрвісу на адвольна зададзеных карыстальнікам уваходных даных. Мноства адпаведнасцяў «уваходныя даныя – выніковыя даныя» можа быць выкарыстана як для аналізу правільнасці працы таго ці іншага элемента ССМТ, так і для стварэння аўтаматызаванага тэставання працы сэрвісу і ССМТ;
- зваротную сувязь ад мэтавай аўдыторыі і экспертаў, у супрацоўніцтве з якімі вядзецца распрацоўка;
- зваротную сувязь ад звычайных зацікаўленых у тэматыцы карыстальнікаў.

Прымаючы да ўвагі пералічаныя вышэй мэты, можна сфармуляваць наступныя асноўныя патрабаванні да праектуемых інтэрнэт-сэрвісаў:

- простасць, зручнасць і інтуітыўная зразумеласць карыстальніцкага інтэрфейсу;
- захаванне ўсіх даных, пададзеных карыстальнікам на ўваход;
- захаванне ўсіх выніковых даных;
- аператыўнае аўтаматызаванае апаўшчэнне распрацоўшчыкаў аб памылках у працы сэрвісаў і сістэмы сінтэзу маўлення.

Праектаванне інтэрнэт-сэрвісаў уключае ў сябе як вызначэнне патрабаванняў, так і апісанне стабільнай структуры і ўзаемадзеяння элементаў сістэмы. У якасці графічнай мовы дакументавання выкарыстоўваецца натацыя UML 2.0 (Unified Modeling Language).

На малюнку 2 прадстаўлена схема кампанентаў узаемадзеяння карыстальніка з шэрагам службовых сэрвісаў сістэмы адкрытых кампанентаў для пабудавання маўленчага інтэрфейсу – www.Corpus.by. Пяройдзем да разгляду функцыянальных магчымасцей праектуемых інтэрнэт-сэрвісаў.

3. Вызначэнне варыянтаў выкарыстання інтэрнэт-сэрвісаў

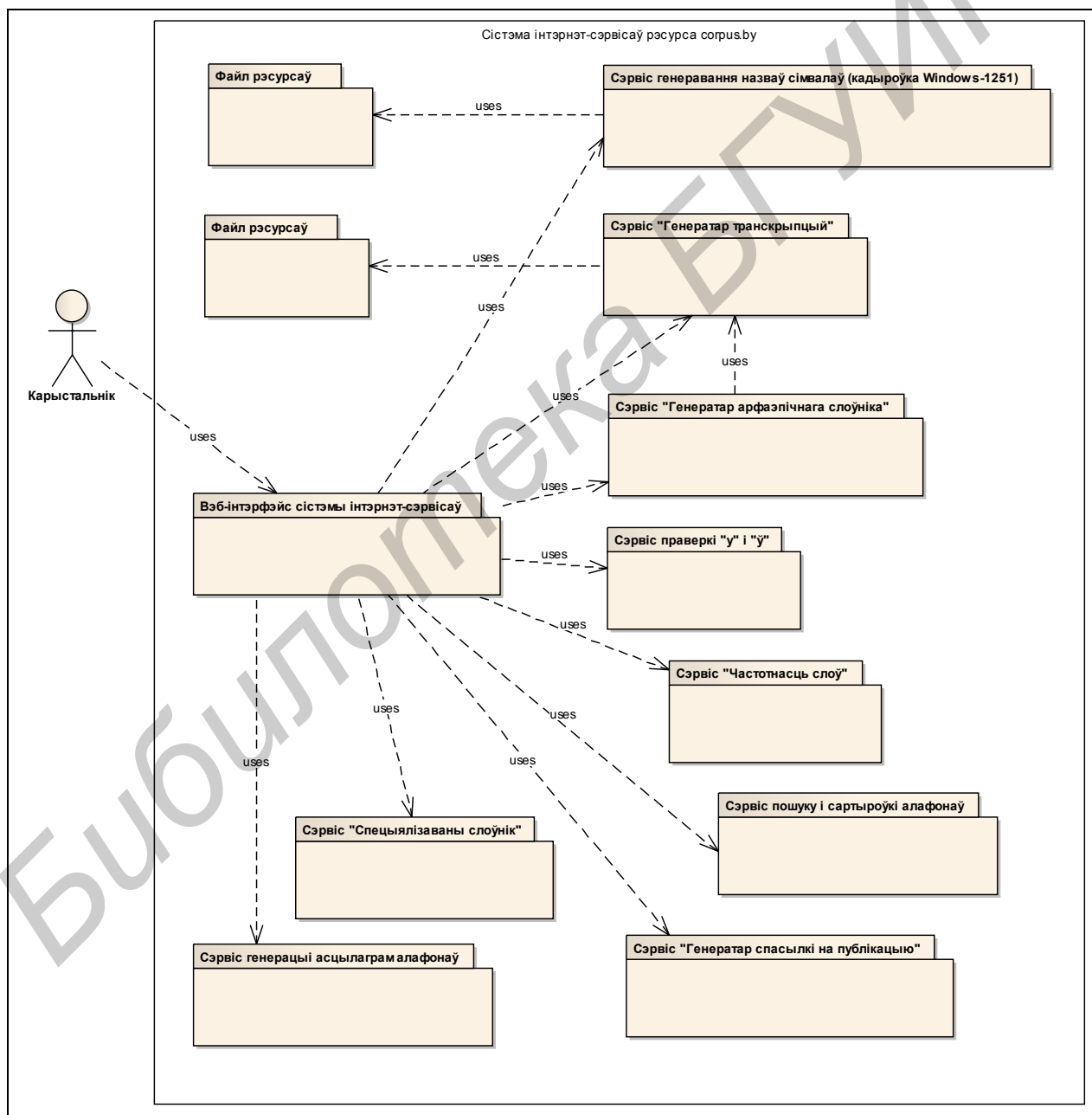
Праектуемая сістэма інтэрнэт-сэрвісаў прадугледжвае выкарыстанне яе як знешнімі

карыстальнікамі, так і распрацоўшчыкамі. Выкарыстанне інтэрнэт-сэрвісаў аднымі і другімі можа як супадаць, так і адрознівацца, бо шэраг распрацоўшчыкаў з'яўляюцца адначасова і карыстальнікамі ў сваіх працоўных мэтах.

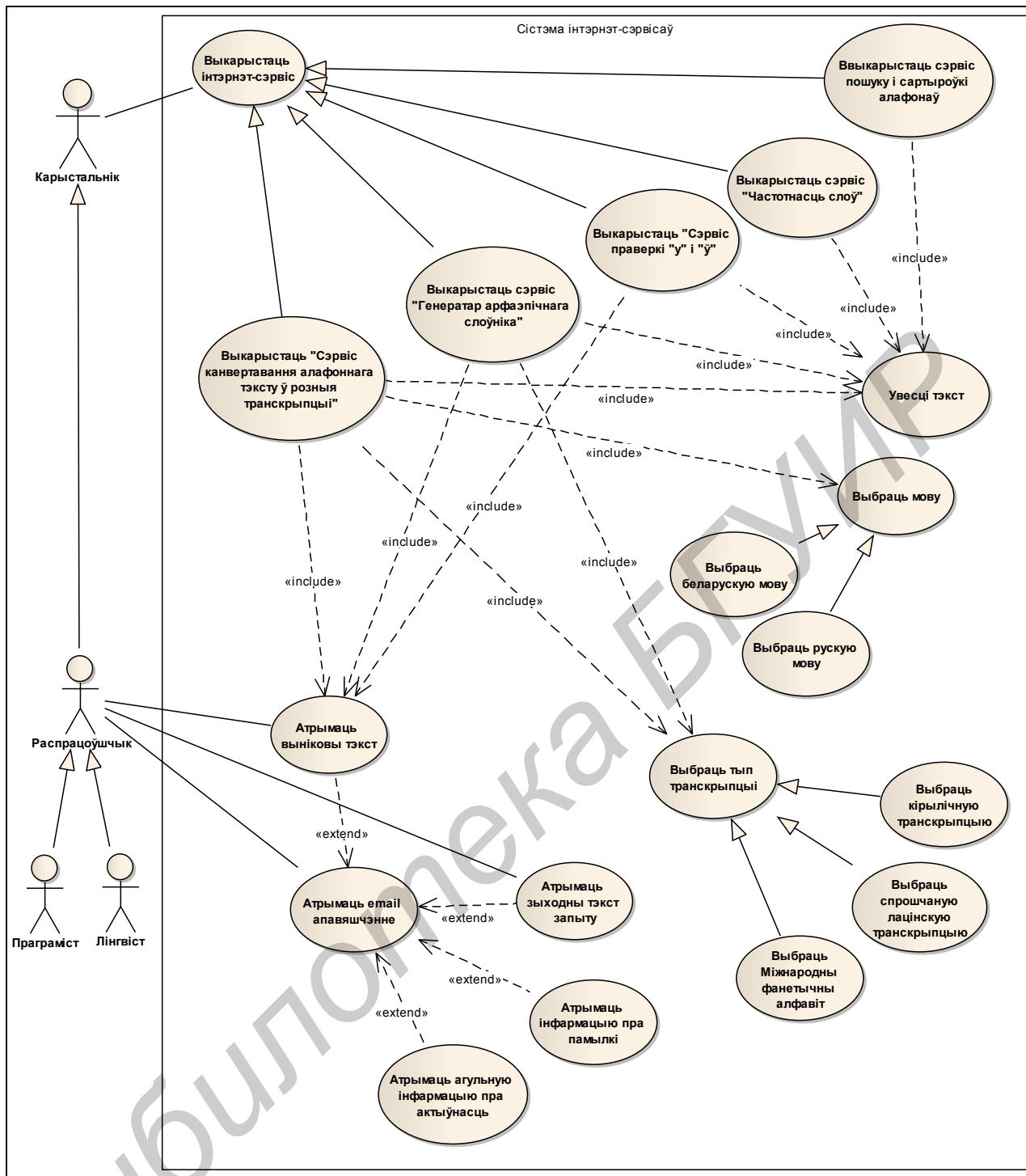
На малюнку 3 прадстаўлена схема варыянтаў выкарыстання сістэмы інтэрнэт-сэрвісаў карыстальнікамі і распрацоўшчыкамі.

На схеме варыянтаў выкарыстання прыведзены наступныя сэрвісы: “Сэрвіс канвертавання алафоннага тэксту ў розныя транскрыпцыі”, “Генератар арфаэпічнага слоўніка”, “Сэрвіс праверкі “y” і “ў”, сэрвіс “Частотнасць слоў”, “Сэрвіс пошуку і сартыроўкі алафонаў”. У кожным

з іх карыстальнік можа падаць на ўваход пэўны тэкст альбо паслядоўнасць сімвалаў і атрымаць на выхадзе неабходны выніковы тэкст у адпаведнасці з задачай, якую вырашае дадзены сэрвіс. У «Сэрвісе канвертавання алафоннага тэксту ў розныя транскрыпцыі» і сэрвісе «Генератар арфаэпічнага слоўніка» карыстальнік можа абраць адзін з трох тыпаў транскрыпцыі: кірылічную, спрошчаную лацінскую ці транскрыпцыю ў фармаце Міжнароднага фанетычнага алфавіту. У «Сэрвісе канвертавання алафоннага тэксту ў розныя транскрыпцыі», як і ў шэрагу іншых, не пазначаных на схеме сэрвісаў, магчымы выбар мовы ўваходнага тэксту: беларускай альбо рускай.



Малюнак 2 – Схема кампанентаў узаемадзеяння карыстальніка з сістэмай інтэрнэт-сэрвісаў www.Corpus.by



Малюнак 3 – Схэма варыянтаў выкарыстання сістэмы інтэрнэт-сэрвісаў www.Corpus.by

Распрацоўшчыкі могуць выкарыстоўваць інтэрнэт-сэрвісы для аналізу працы сэрвісаў і прыняцця рашэнняў па дапрацоўцы і выпраўленні памылак іх працы, а таксама для вырашэння шляхоў развіцця, пошуку сродкаў і спосабаў павышэння эфектыўнасці працы праз назіранне за пэўнымі тэндэнцыямі і заканамернасцямі ў карыстальніцкіх запытах да сэрвісаў.

Дадзеныя сэрвісы праектуюцца як для звычайных карыстальнікаў, так і для адмысловай мэтавай аўдыторыі: лінгвістаў-экспертаў і ўласна

распрацоўшчыкаў праграмных сродкаў: дадзеных сэрвісаў і ССМТ. Такім чынам, распрацоўшчыкі (лінгвісты і праграмісты) з'яўляюцца і прыватным выпадкам карыстальніка, але між іншага валодаюць і сваімі ўласнымі спецыфічнымі магчымасцямі выкарыстання.

Звычайныя карыстальнікі маюць магчымасць бясплатна выкарыстоўваць інтэрнэт-сэрвісы рэсурсу www.Corpus.by у адвольных мэтах. Прывядзём апісанне некалькіх распрацаваных і даступных у інтэрнэце сэрвісаў.

4. Апісанне сэрвісаў: функцыянальныя магчымасці, спосабы выкарыстання

4.1. Генератар файлаў для праслухоўвання

Сэрвіс “Generator of Listen File” ставіць сваёй мэтай прывесці хаця б па адным прыкладзе для адлюстравання кожнага канкрэтнага алафона ці дыфона (для дыфонаў пакуль рэалізаваны толькі выпадак “санорны зычны – галосны”) [Listen file, 2014]. У аснове дадзенага сэрвісу ляжаць спісы слоў, якія былі спецыяльна складзены экспертам. Для зручнасці карыстання і ўспрымання, падабраныя словы згрупіраваны па тэматыках, паколькі адвольны набор слоў, зусім не звязаных між сабой, часта выклікае пэўныя складанасці – карыстальнік падсвядома пачынае суадносіць слова з яго прадметнай вобласцю, што адцягвае ўвагу ад пастаўленай перад ім задачы і запавольвае працу.

Для таго, каб атрымаць спіс слоў для алафонаў, размеркаваных па тэматыках, неабходна націснуць на кнопку “Read Data from Xlsx File (Allophones)” (малюнак 4а), адпаведна для дыфонаў – “Read Data from Xlsx File (Diphones)” (малюнак 4б). Акрамя прыкладаў карыстальнік таксама атрымае інфармацыю пра колькасць тэматык і слоў, а таксама сярэдняю колькасць слоў у адной тэматыцы.

Функцыянальнасць сэрвісу “Generator of Listen File” дае магчымасць паляпшаць працу сінтэзатара маўлення шляхам начытвання адабраных і разбітых па тэматыках слоў і выразання алафонаў, якіх на дадзены момант няма ў базе сінтэзатара і якія маюць ня вельмі добрую якасць. Акрамя таго дадзены сэрвіс можа выкарыстоўвацца для больш хуткага і зручнага стварэння новых “галасоў” для сінтэзатара.

Generator of Listen File

Generator of Listen File

Read Data from Xlsx File (Allophones)

Read Data from Xlsx File (Diphones)

Read Import Txt File

Read Import Xlsx File

Read Data from Xlsx File (Allophones)

Topics = 42, Words = 719, Words per Topic = 17.1

Тэматычны дамен: будоўля..##

аббудава+ць..#
аббі+ць..#
аббіва+нне..#
абпіло+ўванне..#
адту+ліна..#
бу=йнамашта+бных..#
бэ=лечна-кансо+льны..#
гу=каізаля+цыя..#
да+х хі+сткі..#
жыллё+..#
не+калькі а+рак..#
пабудава+ць э+лінг..#

Тэматычны дамен: ваенная тэматыка..##

правядзі+ а=нтывае+нную а+кцыю..#
радыёгра+ма..#
раззбрае+нне..#
ружжо+..#

а)

Generator of Listen File

Generator of Listen File

Read Data from Xlsx File (Allophones)

Read Data from Xlsx File (Diphones)

Read Import Txt File

Read Import Xlsx File

Read Data from Xlsx File (Diphones)

Topics = 37, Words = 805, Words per Topic = 21.8

Тэматычны дамен: абстра+ктыныя назо+ўнікі..##

абплята+нне..#
адда+насць..#
аддале+нне..#
адту+ліна..#
аслупяне+нне..#
бясці+лле..#
выбе+льванне..#
ка+шлянне..#
любо+ў..#
рвань..#
імкне+нне..#

Тэматычны дамен: будо+ўля..##

аб'е+ктавы..#
аб'е+мны..#
аб'е+місты..#
аббудава+ць..#
абпіло+ўванне..#

б)

Малюнак 4 – Прыклад адлюстравання слоў для рэпрэзентацыі а) алафонаў; б) дыфонаў

4.2. Генератар спасылкі на публікацыю

Стварэнне бібліяграфічных спасылак – неад’емная частка навуковай працы. Каб спрасціць гэтую задачу для карыстальнікаў, быў створаны сэрвіс “Get a publication reference!” [Publication references, 2014]. Дадзены сэрвіс распрацаваны для аўтаматычнай генерацыі спасылак.

Зайшоўшы на старонку сэрвісу, карыстальніку неабходна пазначыць фармат спасылкі, тып крыніцы (кніга, артыкул часопісу, канферэнцыя, вэб-сайт) і мову афармлення публікацыі (беларуская, руская ці англійская) (малюнак 6).

Наступным крокам з’яўляецца ўвод даных карыстальнікам. Абавязковыя для запаўнення палі пазначаны чырвонай зорчак (*). Таксама прадугледжана магчымасць прагляду прыкладаў запаўнення, для чаго неабходна націснуць на кнопку “Ачысціць усё і паказаць прыклады”. Прыклад запаўнення формы данымі паказаны на малюнку 7.

Малюнак 6 – Пазначэнне фармату, тыпу і мовы публікацыі

Малюнак 7 – Запаўненне формы данымі

Па заканчэнні запаўнення формы неабходна націснуць на кнопку “Атрымаць спасылку на публікацыю”. Вынік з’явіцца ў асобным полі ніжэй (малюнак 8).

На дадзены момант для карыстальнікаў даступны толькі ВАК-фармат публікацый. Плануецца пашырыць магчымасці сэрвісу праз павелічэнне колькасці тыпаў спасылак і далучэнне наступных фарматаў: Chicago, MLA, APA.

Publication references

Гецэвіч, Ю.С. Стварэнне сэрвіса арфаэпічнага генератара слоўнікаў / Ю.С. Гецэвіч, А.В. Гюнтар, С.І. Лысы [і інш.] // Тэзі доповідей міжнародной конференції «Діалекты в синхроніі та дяхроніі : загальнаслов’янський контекст» (Кіев, 2–4 квітня 2014 року) / За ред. П.Ю. Гриценка. Ін-т укр. мови НАН Украіны. Кіев : КММ, 2014. — С. 101-106.

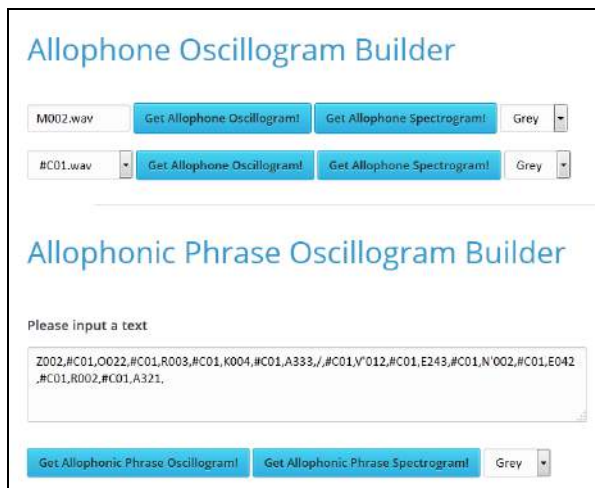
Малюнак 8 – Згенераваная спасылка на публікацыю

Праца з генератарам бібліяграфічных спасылак дазволіла даведацца больш дакладна пра іх змест і структуру, што з’яўляецца добрай асновай для распрацоўкі сістэмы распазнавання спасылак у тэкстах беларускай і рускай моў. Інфармацыя якую можна атрымаць з бібліяграфічных спасылак пры распазнаванні, можа быць карыснай для бібліятэк, выдавецтваў ці навуковых устаноў, бо яна змяшчае звесткі пра аўтараў і рэдактараў, назвы публікацый, назву і месца выдавецтва, год выдання і інш.

4.3. Сэрвіс графічнага адлюстравання алафонаў і алафонных фраз

Сэрвіс графічнага адлюстравання алафонаў і алафонных фраз “Allophone Plotter” – гэта яшчэ адзін са службовых сэрвісаў сістэмы адкрытых кампанентаў для пабудавання маўленчага інтэрфейсу – www.Corpus.by [Allophone Plotter, 2014]. Яго галоўнай задачай з’яўляецца графічнае адлюстраванне фізічнага сігналу ў часавым альбо частотным выглядзе адпаведнай фанетычнай адзінкі – алафона, абранага карыстальнікам сэрвісу, альбо платаванне цэлай фразы па фанетычным тэксце, які мае форму радка з паслядоўна запісаных алафонаў. Для адлюстравання сігналу ў часовай вобласці выкарыстоўваецца графік асцылаграмы – залежнасці амплітуды сігналу ад часу. Для частотнага прадстаўлення выкарыстоўваецца графік спектраграмы – залежнасці амплітуды ці энергіі сігналу ад часу і ад частаты адначасова. У гэтым выпадку графік выглядае як двухмерны каляровы малюнак сігналу з часам, які адкладаецца па восі абсцыс, частатой – па восі ардынат, а ўзровень энергіі сігналу характарызуецца інтэнсіўнасцю колеру на малюнку.

Сэрвіс можа выкарыстоўвацца экспертамі-фанетыстамі, лінгвістамі, студэнтамі філалагічных і педагогічных вузаў ці проста зацікаўленымі асобамі для знаёмства з “выглядам” і фізічнымі характарыстыкамі алафонаў (рэалізацыямі фанем), у працэсе вывучэння беларускай мовы. Напрыклад, спектраграма можа вельмі моцна дапамагчы ў вывучэнні як фанетыкі мовы ў агульным, так і асобных гукаў мовы ў прыватнасці. Кожны алафон мае свае асаблівасці з пункту гледжання яго фізічных параметраў, якія з’яўляюцца вынікамі дзеяння інтра- і экстралінгвістычных фактараў мовы і якія складаюць “жывую непаўторную карціну” гэтага гука, добра бачную на частотна-часавым плане, які сабой уяўляе спектраграма.

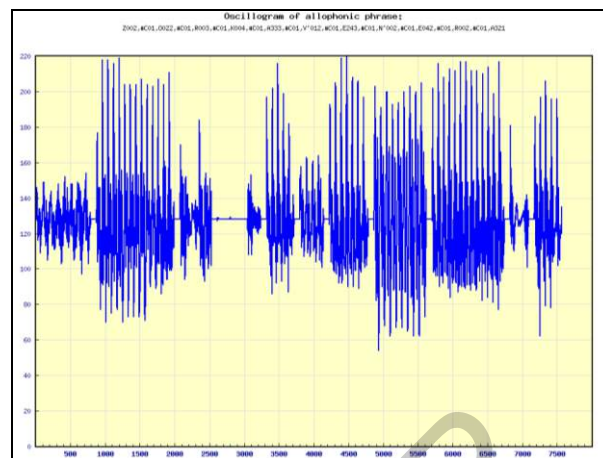


Малюнак 9 – Знешні інтэрфейс “Сэрвісу графічнага адлюстравання алафонаў і алафонных фраз”

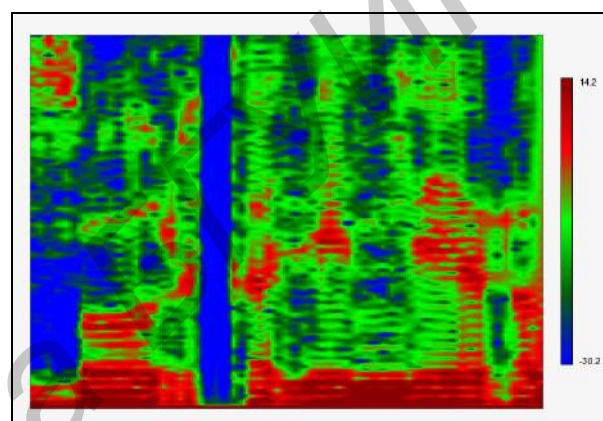
Асцылаграма фанетычнай фразы можа эфектыўна выкарыстоўвацца для вывучэння прасодыкі (інтанацыйнага склада) мовы, ад якой вельмі моцна залежыць успрыманне маўлення чалавекам. Як вядома, прасодыка складаецца з трох кампанентаў: энергетыкі – цяжучай змены сілы гуку, рытмікі – змены працягласці гукаў і паўз і мелодыкі – руху частаты асноўнага тону. Усе гэтыя моўныя з’явы яркава назіраюцца на асцылаграме. Напрыклад, праз адлюстраванне рытмікі асцылаграма можа падказаць карыстальніку, дзе адносна агульнай працягласці фразы, знаходзіцца слоўны ці фразавы націск, таму што ў гэтым месцы на графіку павінна назірацца большая амплітуда сігналу.

На сённяшні момант агульны сцэнар працы з сэрвісам выглядае наступным чынам. Карыстальнік заходзіць на старонку “Сэрвісу графічнага адлюстравання алафонаў” (знешні інтэрфейс старонкі прадстаўлены на малюнку 9). На гэтай старонцы ён можа наплатваць асобны алафон ці алафонную фразу. Асобны алафон можна ўвесці ўручную ў тэкставым полі ці выбраць з фанетычнай базы праз выпадаючае меню. Алафонную фразу карыстальнік павінен уводзіць у спецыяльна адведзенае тэкставае поле, у якім дазволены шматрадковы ўвод. Пасля ўвода асобнага алафона ці алафоннага тэксту карыстальнік павінен націснуць адпаведную кнопку для пабудавання асцылаграмы ці спектраграмы. Пры выбары платавання спектраграмы, у карыстальніка ёсць магчымасць выбару карты колеру: чорна-белай ці каляровай. Далей трэба пачакаць некаторы час пакуль сэрвіс апрацуе запыт і згенеруе адказ у выглядзе html-старонкі, у якую ўбудаваны малюнак сігналу ў фармаце png, і якую браўзер адлюструе карыстальніку. Вынікі працы сэрвісу прадстаўлены на малюнку 10.

Асаблівацю сэрвісу з’яўляецца магчымасць апрацоўкі сігналу і будавання графічных “партрэтаў гукаў жывой мовы on-line”, а таксама яго арыентацыя на беларускую мову.



а)



б)

Малюнак 10 – Вынікі працы “Сэрвісу графічнага адлюстравання алафонаў і алафонных фраз”: а) асцылаграма алафоннай фразы; б) спектраграма алафоннай фразы

На дадзены момант аўтары не маюць інфармацыі аб аналагічных бясплатных on-line сэрвісах для лінгвістаў ці людзей, зацікаўленых у вывучэнні беларускай мовы, якія валодалі б такой жа самай функцыянальнасцю. У далейшым плануецца выкарыстанне гэтага сэрвісу ў якасці дадатковай функцыянальнасці для больш высокаўзроўневых сэрвісаў рэсурсу www.Corpus.by.

Заклучэнне

У выніку праведзенай працы былі распрацаваны агульныя падыходы да стварэння сістэмы інтэрнэт-сэрвісаў www.Corpus.by у мэтах паляпшэння і ўдасканалення працы сістэмы сінтэзу маўлення па тэксце з адначасовай магчымасцю як вырашэння розных прыкладных задач, так і інтэграцыі з іншымі праектамі, у прыватнасці з праектам OSTIS. У дадзеным артыкуле апісаны асаблівасці праектавання, рэалізацыі і выкарыстання сістэмы, прыведзены апісанні некаторых распрацаваных сэрвісаў, а менавіта: сэрвісу “Генератар файлаў для праслухоўвання”, сэрвісу “Генератар спасылкі на публікацыю”, “Сэрвісу графічнага адлюстравання алафонаў і алафонных фраз”.

Бібліяграфічны спіс

[Allophone Plotter, 2014] Allophone Plotter [Electronic resource]. – 2014. – Mode of access : <http://corpus.by/readVoiceDBAllophones/>. – Date of access : 12.12.2014.

[Listen file, 2014] Generator of Listen File [Electronic resource]. – 2014. – Mode of access : <http://corpus.by/genListenFile/>. – Date of access : 12.12.2014.

[Publication references, 2014] Get a publication reference! [Electronic resource]. – 2014. – Mode of access : <http://corpus.by/publicationReference/>. – Date of access : 12.12.2014.

[www.Corpus.by, 2012] Text-to-Speech PHP-Based Synthesizer [Electronic resource]. – 2012. – Mode of access : <http://corpus.by/>. – Date of access : 12.12.2014.

[Гецэвіч, 2012] Гецэвіч, Ю.С. Алгарытмы лінгвістычнай апрацоўкі тэкстаў для сінтэзу маўлен-ня на беларускай і рускай мовах: дыс. ... канд. тэхн. навук / Ю.С. Гецэвіч. – Мінск, 2012. – 193 с.

[Лобанов, 2008] Лобанов, Б.М. Компьютерный синтез и клонирование речи / Б.М. Лобанов, Л.И. Цирульник. – Минск : Белорусская наука, 2008. – 344 с.

WWW.CORPUS.BY: OPEN-SOURCE COMPONENTS FOR NATURAL LANGUAGE INTERFACES

Hetsevich Y.S.*, Lobanov B.M.*, Lysy S.I.*,
Hiuntar E.V.*, Denisyuk D.A.*,
Zakharyeu V.A.**

* *United Institute of Informatics Problems,
National Academy of Sciences, Minsk, Belarus*

**yury.hetsevich@gmail.com,
lobanov@newman.bas-net.by,
stanislau.lysy@gmail.com, lena205593@gmail.com,
d.denisyuk@gmail.com**

** *Belarusian State University of Informatics and
Radioelectronics, Minsk, Belarus*
delfvad@gmail.com

This paper describes general approaches to designing, implementation and usage of open-source components of www.Corpus.by internet-resource for Natural Language Interfaces. These components serve to improve the performance of a text-to-speech system (TTS), while at the same time providing the opportunity of integration with the OSTIS project. The paper also define goals and objectives, ways of operation and available functionality of some services that have been implemented.

INTRODUCTION

In order to solve the problem of speech interaction between people and machines, scientists all over the world are working on creating speech synthesis and recognition systems. Recently, TTS systems have reached certain degree of accuracy and become widely used, but it is impossible to say that the task of text-to-speech synthesis is completely resolved, especially it concerns the Belarusian language. In this regard, the authors decided to create an Internet resource that would bring together a number of Internet-services,

each aimed at solving certain tasks and at the same time helping developers to make the TTS system performance more optimal and precise.

MAIN PART

The tasks facing the developers of a text-to-speech system may be conventionally divided into three groups: preprocessing of input texts, speech synthesis itself and processing of resulting data.

The internet-resource www.Corpus.by is aimed at addressing the challenge of text-to-speech synthesis itself, as well as at solving a range of additional tasks of text-preprocessing and processing of the results of the text-to-speech synthesizer's operation. In order to achieve this aim, the decision was made to develop a system of Internet-services, all of each would solve its own certain task and, when combined with all the rest services, would also solve the essential subtasks of the text-to-speech system. The following services are currently developed or are on the intermediate stage of development and are easily available on the Internet:

- “Service of searching and sorting of allophones”;
- “Service of converting allophonic texts into different transcriptions”;
- “Service of generating names of characters (encoding Windows-1251)”;
- “Transcription Generator”;
- “Orthoepic Dictionary Generator”;
- “Spell-check for Belarusian U”;
- “Frequency of Words”;
- “Get Publication References”;
- “Allophone Oscillogram Builder”;
- “Allophone Plotter”;
- “Get Listen File”;
- “Sound Recorder”;
- “Specialized Dictionary”;
- etc.

CONCLUSION

As a result of the work carried out by the authors, in order to improve and make more precise the performance of the text-to-speech system, and at the same time to provide free electronic services to the population in order to resolve a wide range of applied tasks, general approaches to creation of the system of internet-services www.Corpus.by were developed.