

Недостатком является возможность возникновения ошибки при формировании словоформ. Поэтому дальнейшие оптимизации и модификации данного алгоритма являются актуальными.

Список литературы

1. Грамматика кыргызского языка: краткий справочник для студентов, Бишкек 2002
2. Глазунова О. Грамматика русского языка в упражнениях и комментариях (Морфология + Синтаксис).
3. Θ. Калыева Разговорник русско-кыргызский - Орусчакыргызча сүйлөшмө
4. Исаилова Н.А. Принципы организации морфологического анализатора в трансляторе. /Известия КГТУ им. И. Рazzакова-2010, №20, С233-236
5. Исаилова Н.А. Организация морфологического анализа в трансляторах. /Вестник Восточно-Казахстанского государственного технического университета им. Д. Серикбаев, Усть-Каменогорск, №1, март, 2012 г.-С 97-101
6. Исаилова Н.А. Алгоритмы отладки процесса трансляции. /Известия КГТУ им. И. Рazzакова-2011, №22, С278-279

УДК 004.432.4

AUTOMATING OF THE TEXT GENERATION WITH A GIVEN REPRESENTATIVENESS PHONETIC UNITS BY A FORMALIZATION OF PHONOLOGICAL RULES

Aigerim Buribayeva, PhD, L.N.Gumilyov Eurasian National University, 010008, Astana, Pushkin str. 2, e-mail: buribayeva@mail.ru

Arman Kaliyev, L.N.Gumilyov Eurasian National University, 010008, Astana, Pushkin str. 2, e-mail: kaliyev.arman@yandex.kz

Banu Yergesh, L.N.Gumilyov Eurasian National University, 010008, Astana, Pushkin str. 2, e-mail: saturn_banu@mail.ru

The article describes formalization of phonological rules of Kazakh language and use it to automate the process of formation carried the text of the material with a given phonetic units of representativeness (in particular diphones). This is necessary, particularly in the development of automatic speech synthesis. Using this versatile program, it will be able to get the text material with full coverage of all possible diphones for all Turkic language.

Keywords: Speech synthesis, text analyzer, sound units, diphones, statistics, acoustic database, text body.

АВТОМАТИЗАЦИЯ ФОРМИРОВАНИЯ ТЕКСТОВОГО МАТЕРИАЛА С ЗАДАННОЙ ПРЕДСТАВИТЕЛЬНОСТЬЮ ФОНЕТИЧЕСКИХ ЕДИНИЦ С ПОМОЩЬЮ ФОРМАЛИЗАЦИИ ФОНОЛОГИЧЕСКИХ ПРАВИЛ

Бурибаева Айгерим Кеулимжановна, PhD, ЕНУ им. Л.Н. Гумилева, 010008, Астана, Пушкина 2, e-mail: buribayeva@mail.ru

Калиев Арман Куанышевич, ЕНУ им. Л.Н. Гумилева, 010008, Астана, Пушкина 2, e-mail: kaliyev.arman@yandex.kz

Ергеш Бану Жантуганкызы, ЕНУ им. Л.Н. Гумилева, 010008, Астана, Пушкина 2, e-mail: saturn_banu@mail.ru

В статье описана формализация фонологических правил казахского языка и на ее основе осуществлена автоматизация процесса формирования текстового материала с заданной представительностью фонетических единиц (в частности дифонов). Это

необходимо, в частности, при разработке систем автоматического синтеза речи. Так, используя данную универсальную программу, появится возможность получить текстовый материал с полным покрытием всех возможных дифонов для любого тюркского языка.

Ключевые слова: Синтез речи, текстовый анализатор, звуковые единицы, дифоны, статистика, акустическая база, текстовый корпус

Introduction: Many modern trends of development speech technologies involve the use of speech databases. These databases are based on the texts in the utterance of one or more speakers. The criterion of suitability speech serves as the base, above all, the fullness its speech elements. For example, for the development of speech synthesis systems, such elements may be different depending on the selected base units. Most often it is diphones or allophones. For the diphone's synthesis requires a database, containing all possible for a given language binomial combination of phonemes (allophones).

In this connection with the above, particular importance is the preparation of text material in problems where a set of necessary elements of speech pre-defined framework, especially in cases where the resource for processing and structuring of the speech material is limited.

In addition to the phonetic representation, using a special text material provides a fixed amount of compactness. This further allows to significantly reduce the time to process it. It is believed that the use of large amounts of textual material completeness of the coverage units is achieved without special selection, however, this approach inevitably arises redundant database (which may in the future to require post-processing of the material to exclude duplicates). Also, some rare phonetic combinations that occur at the junction of words, may not once meet. In order to simplify and automate the process of formation of the texts with a given phonetic representativeness, assessing the completeness of coverage of items in a given text, as well as reducing the body text of redundancy due to the removal of a recurring items, was developed by the phonetic analyzer text material.

A feature of this program is its versatility. It can be used for virtually any Turkic language, even without the formalization of phonological rules, it is enough just to have a complete pronouncing dictionary rather than spelling.

The formalization of the phonological rules of sound combinations in Kazakh language: For the development of phonetic transcriptor were investigated orthoepic rules of the Kazakh language. For the convenience of the reader in this text, the rules are divided into groups that are numbered:

1. In the Kazakh language, if the word begins with vowel «е», then in front of it pronunciation is heard as «й», if the word starts with the vowel «о», «ө», then the pronunciation in front of them formed a brief insert «ү» for example, «ет» – «йет», «он» – «уюн», «өнер» – «үөнер».

2. If a word begins with a consonant «р» or «л», then the pronunciation of these sounds could be heard before the vowel «ы», «и», depending on the hardness or softness of consonants here «р», «л» means soft analogs "«п» and «л»". For example, «рас» – «ырас», «рет» – «ирет», «лас» – «ылас», «лездे» – «ілездे».

3. When pronouncing borrowed sound «ю» as part of word is heard «йүү», «йүү», depending on the hardness or softness of the other vowels in syllables. For example: «қою» – «қойүү», «тую» - «түйүү»;

4. When pronouncing borrowed sound «я» as a part of a word is heard «йа», «йә», depending on the hardness or softness of the other vowels in syllables. For example: «аян» – «айан», «элия» – «әлійә»;

5. When pronouncing borrowed sound «и» in a consisting of words heard «ый», «ий», depending on the hardness or softness of the other vowels in syllables. For example, «ине» – «ійне», «жина» – «жыйна». If before or after «и» go according to «к», «ғ» with descender, that the pronunciation of the sound «и» always heard «ый». For example, «қын» – «қыйын», «қигаш» – «қыйғаш».

6. When the pronunciation of the diphthong «ү» as a part of a word is heard «ұү», «үү», depending on the hardness or softness of the other vowels in syllables. For example, «туыс» – «тұуыс», «куту» – «кутүү».

7. The Vowels «ұ», «ү», «о», «ө» at the beginning or the first syllable of the word in the pronunciation change in the next syllable vowel sounds «ы», «и» on the vowels «ұ», «ү» respectively. For example, «қолтық» – «қолтұқ», «құлын» – «құлұн», «құлкі» – «құлқұ», «қөлік» – «қөлүк»;

8. Vowels «ү», «ө» in the beginning or in the first syllable of the word during the pronunciation changes in the following syllables vowel «е» to the next vowel «ө», for example, «ұлкен» – «ұлкөн», «өнер» - «өнөр».

9. Vowels «ә», «ү», «и» in the beginning or in the first syllable of the word during the pronunciation changes in the following syllables vowel «а» to its allophone «ә», for example, «ләззат» – «ләззәт», «діндар» – «діндәр».

10. If in a word sounds «с» and «ш», «с» and «ж» or «з» and «ш» meet in succession, then instead of them is pronounced the sound of double «шш». Also, instead of borrowed sound «щ» is pronounced «шш». For example: «досжан» - «дошшан», «басшы - башшы», «сөзшең – сөшшөң», «көшсен»-«көшшөң», «аңы»-«ашшы».

11. If in a word after sounds «з» and «ж» meet in succession, instead of them pronounce dual sound «жж», if sounds «з» and «с» meet in succession, instead of them pronounce dual sound «сс». For example: «бозжорға» - «божжорға», «азсыну - ассынұу».

12. If in a word after sound «н» встречается «б» или «п», then the pronunciation of sound «н» replaced to «м». For example: «мінбер – мімбер», «ойынпаз» – «ойымпаз».

13. If in a word after sound «н» meet «г», «ғ», «к» or «қ» then pronounce of «н» replaced to «ң». For example: «тұнгі» – «тұңғұ», «қашанғы» - «қашанқы», «зиянкес» – «зыйанкес», «сәңкөй» – «сәңқөй».

14. When pronouncing the word in the composition of sound combinations әл, ән, әл between two sounds is formed a brief insertion of vowels «ы», «и», depending on the hardness and softness corresponding syllable. For example, «мемлекет» – «мемілекет», «бағлан» – «бағылан», «яғни» – «йағыный».

15. Uncombinable sounds found in many compound words are replaced by the pronunciation sound. For example, «шашбай» - «шашпау», «атбегі»-«атпегі», «атжалман» - «атшалман», «Көбосын» - «Көппосұн», «түпдерек» - «тұбдөрөк», «көпжиын» - «көбжыйын», «көпмүшө - «көбмүшө», «тұпнегіз» - «тұбнегіз», «тасбауыр» - «таспауұр» и т.д.

Transkriptor implemented as a program that replaces some other characters in accordance with the rules contained in the control file. Rules are written in accordance with the each item of orthoepic rules of Kazakh language:

1. #e=ье, #o=yo, #ө=уө;

2. #л^a=ыл^a, #л^o=ыл^o, #л^ү=ыл^ү, #л^ә=іл^ә, #л^ү=іл^ү, #л^e=іл^e, #л^i=іл^i, p^a=ыр^a, #p^o=ыр^o, #p^ү=ыр^ү, #p^ә=ір^ә, #p^ү=ір^ү, #p^e=ір^e, #p^i=ір^i;

3. аю=айұу, ою=ойұу, үю=үйұу, ыю=ыйұу, үю=үйұу, ею=ейұу, қиу=қыйұу, #тию#=тійұу, қиу=кійұу, #сию#=сыйұу, #жию#=жыйұу, а^ио=a^ыйұу, о^ио=o^үйұу, ү^ио=ү^үйұу, ы^ио=ы^ыйұу, ә^ио=ә^ійұу, ө^ио=ө^үйұу, ү^ио=ү^үйұу, і^ио=i^ійұу, е^ио=e^ійұу;

4. ая=айа, оя=ойа, үя=үйа, ыя=ыйа, қия=қыйа, #сия=сыйа, #жия=жыйа, #мия=мыйа, #зия=зыйа, а^ия=a^ыйа, о^ия=o^үйа, ү^ия=ү^үйа, ы^ия=ы^ыйа, ә^ия=ә^ійә, ү^ия=ү^үйә, ия^a=ыйа^a;

5. #ми=мый, #жи=жый, а^и=a^ый, о^и=o^ый, ү^и=ү^ый, ы^и=ы^ый, ә^и=ә^ій, ө^и=ө^ій, ү^и=ү^ій, і^и=i^ій, е^и=e^ій, и^a=ый^a, и^o=ый^o, и^ү=ый^ү, и^ы=ый^ы, и^ә=ій^ә, и^ө=ій^ө, и^ү=ій^ү, и^и=i^ій, и^e=ій^e, қи=қый, ғи=ғый, иқ=ыйқ, иғ=ыйғ;

6. а^у=а^ყу, о^у=օ^ყу, ყ^у=ყ^ყу, ы^у=ы^ყу, ə^у=ə^ყу, ө^у=ө^ყу, ү^у=ү^ყу,
і^у=i^ყу, e^у=e^ყу, y^a=ყу^a, y^o=ყу^o, y^ყ=ყу^ყ, y^ы=ყу^ы, y^ә=ყу^ә, y^ө=ყу^ө, y^ү=ყу^ү,
y^і=ყу^і, y^e=ყу^e;
 7. о^ы=o^ყ, ყ^ы=ყ^ყ, ө^і=ө^ყ, ү^і=ү^ყ;
 8. ө^е=ө^ө, ү^е=ү^ө;
 9. і^а=i^ә, ү^а=ү^ә, ө^а=ө^ә;
 10. сш=шш, сж=шш, зш=шш, шс=шш, ш=шш;
 11. зж=жж, зс=сс;
 12. нб=мб, нп=мп;
 13. нг=нг, нғ=нғ, нк=нқ нқ=нқ;
 14. мл = міл, ғн=ғын, ғл=ғыл;
 15. шб=шп, тб=тп, тж=тш, пб=пп, пд=бд, пж=бж, пм=бм, пн=бн, сб=сп, сд=ст, қб=қп,
қг=қг, қд=ғд, қж=ғж, қз=ғз, қм=ғым, қн=ғын, қб=қп, қг=ққ, қд=қт, қж=ғж, қз=ғз, қм=ғм,
қн=ғн, զկ=զգ, զլ=զբ, զտ=ստ, կլ=ғыլ.

Each substitution rule is composed of two parts separated by a sign «=>. From the left of this sign are original alphabetic character for word recording, on the right - the characters that should be replaced in the transcription.

For transcription of the given word consistently searches next occurrence of the left part of rule in it. If any of it detected, then instead of it, inserted the right part of the rule.

As a transcription symbols for vowels used mainly relevant Kazakh letters. Solid consonants are transcribed as Kazakh letters and the relevant soft consonants with the analogous Latin letters.

«#» means the beginning or end of word, depending on the location of: if «#» standing in front of characters then it is the beginning of word; if «#» stands after the characters it's the end

«^» means any characters in any number between two sounds

Each substitution rule is composed of two parts separated by a sign «=>. From the left of this sign are original alphabetic character recording word on the right - the characters that should be replaced in the transcription.

For transcription of given word consistently searches next occurrence of the left part of rule in it. If any of it is detected, then it is inserted along the right side of the rule.

It is recommended that in the control file of these groups in numerical order, without changing the order of the rules in groups because the order of substitutions is obviously important.

Also orthoepic rules were included to transcriptor rules defining the softness and labial consonants:

2. $\theta\bar{b}=eb$, $\theta\bar{g}=eg$, $\theta\bar{d}=ed$, $\theta\bar{v}=ev$, $\theta\bar{z}=ez$, $\theta\bar{j}=ej$, $\theta\bar{k}=ek$, $\theta\bar{l}=el$, $\theta\bar{m}=em$, $\theta\bar{n}=en$, $\theta\bar{q}=eq$, $\theta\bar{o}=ef$, $\theta\bar{r}=er$, $\theta\bar{s}=es$, $\theta\bar{t}=et$, $\theta\bar{u}=eu$, $\theta\bar{w}=ew$, $\gamma\bar{b}=yb$, $\gamma\bar{g}=yg$, $\gamma\bar{d}=yd$, $\gamma\bar{v}=vv$, $\gamma\bar{z}=yz$, $\gamma\bar{j}=yj$, $\gamma\bar{k}=yk$, $\gamma\bar{l}=yl$, $\gamma\bar{m}=ym$, $\gamma\bar{n}=yn$, $\gamma\bar{h}=yh$, $\gamma\bar{f}=yf$, $\gamma\bar{p}=yr$, $\gamma\bar{c}=ys$, $\gamma\bar{t}=yt$, $\gamma\bar{u}=yu$, $\gamma\bar{w}=yw$, $\delta\bar{e}=be$, $\theta\bar{e}=ge$, $d\bar{e}=de$, $\bar{e}\theta=ve$, $z\bar{e}=ze$, $\bar{e}j=j\theta$, $k\bar{e}=ke$, $l\bar{e}=le$, $m\bar{e}=me$, $n\bar{e}=ne$, $q\bar{e}=qe$, $f\bar{e}=fe$, $r\bar{e}=re$, $s\bar{e}=se$, $t\bar{e}=te$, $u\bar{e}=ue$, $w\bar{e}=we$, $b\bar{y}=by$, $g\bar{y}=gy$, $d\bar{y}=dy$, $v\bar{y}=vy$, $z\bar{y}=zy$, $j\bar{y}=jy$, $k\bar{y}=ky$, $l\bar{y}=ly$, $m\bar{y}=my$, $n\bar{y}=ny$, $q\bar{y}=qy$, $p\bar{y}=fy$, $r\bar{y}=ry$, $s\bar{y}=sy$, $t\bar{y}=ty$, $u\bar{y}=uy$, $w\bar{y}=wy$;

Let us explain marks used in the replacement rules. Latin characters in a group of 16 means that the sound is soft and non-labial, 17 in group "2" after the consonant means that the consonant - solid lip and in the group of 18 the number "3" after the consonant means that the consonant - soft [labial]

In general phonetic transcriptor was about 400 rules. Similar rules have been further used for the transcription of sounds borrowed from the Russian language, formalized in [1].

Description of the program: The program involves automating the process of forming the texts of the material with a given phonetic units of representativeness (in particular diphones). This is necessary for development of automatic speech synthesis. Thus, using this program, you will be able to get a text material with full coverage of all possible diphones in Kazakh language.

Of course, we are not talking about the original generation of coherent text – we cannot afford this machine. This refers to the layout of the text matching words contained in pronouncing dictionary. Selection of the program will be implemented in such a way as to minimize redundancy (repetition) of units in the text.

Thus, the program allows the following tasks:

- the generation of the primary list of diphones submitted by Alphabet;
- checking for the specified diphones in the pronouncing dictionary, and a definitive list of the correct diphones;
- generation of mini-set of the text covering all sorts of diphones;
- evaluation of the degree of phonetic informativeness of the words included in the text material.

As mentioned above, the program can be used for virtually any Turkic language, without the formalization of phonological rules, replacing the spelling dictionary to the pronouncing dictionary. For Kazakh language, we used the spelling dictionary of 45,000 word forms.

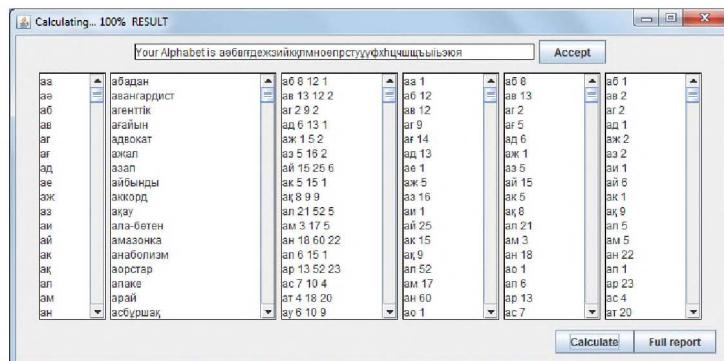


Fig. 2. Text analyzer's window

The first field is displayed diphones occurring at the beginning, in the middle and at the end of the words in the second - found only at the beginning of the word in the third - found only in the middle of the word in the fourth - found only at the end of the words, with appropriate statistics.

By clicking on the "Full report" we get a detailed report with a set of text.

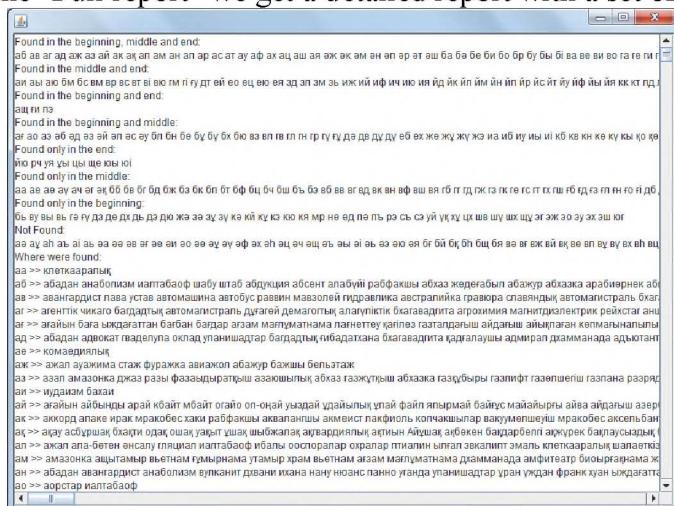


Fig. 3. Full report's window

A selection of the text carried by the following algorithm:

1. Get over the words for a particular diphones;
 2. The words are divided into diphones;
 3. Each diphone in words is checked by a priority (the number of occurrences of words chosen before);
 4. Selection of the words with the lowest priority;
 5. New priorities are assigned to each diphone in the selected word.

Of course, there are common diphones, for example, diphone «ah» in the text set occurs 100 times, when diphone «ae» occurs only 1 time. Yet, due to the method of assigning priorities to phonetic units, we obtain a set of text less redundant and compact.

Results:

Quantity of diphones in the primary list 412=1681;

After checking on existence in their pronouncing dictionary remained – 1003:

From them:

Meeting at the beginning, in the middle and at the end of the word – 297:

Meeting at the beginning and in the middle of the word – 144:

Meeting at the beginning and in the middle of the word
Meeting at the beginning and at the end of the word – 3:

Meeting in the middle and at the end of the word – 159.

Meeting only at the beginning of the word - 45

Meeting only in the middle of the word - 347

Meeting only in the middle of the word - 57;
Meeting only at the end of the word - 8:

Quantity of words in a text set = 6131

Quantity of words in a text set = 6151.
Found at the beginning, in the middle and at the end.

аб ав аг ад аж аз ай ак ак ад ам ан ад а

ао ав аі ад аж аз аи ак ақ ал ам ан аи ар ас аі ау ае
бә би ба бр бу бүл біл вә ве ви во гә ге гү го гүл гәл

ба бе би бо бу бы ог ва ве ви во га гети то ту та ғы да дә дж ди до др ду ді ев ег ед еж
ез ек ел ем ен еп ер ес ет еу еф еш жа же жи жу жы жі за зе зи зо зу зы зі ив ид ие из ии
ик ил им ин ио ип ир ис ит их иш йе йо йі ка ке ки кл ко кр кс ку кх кі қа қи қу қы ла ле ли лл
ло лу лы ль ма мб ме ми мн мо му мы мі на нә не ни но ну ны ні ню об ов ог од ож оз ой ок
ок ол ом он оо оп оп ос от оф ош ою өз өк өл өн өс па пе пи по пр пс пт пу пы пі ра рә ре ри
ро ру ры са се си ск см со сс ст су си сі сь та те ти то тр ту ты ті ть уа уг уе уз уи ук ул ум ун
уп ур ус ут уш уы уі үз үк үл үм үн үп үр үс үя үз үй үк үл үн үр үс үт үш үю фа фо фр фт
фь ха ца ча чи ша ше ши шт шу шы ші ыб ыз ық ыл ым ын ыр ыс ыт ыш іш іл ін ір іс іш эз эл
эн эр эт юк юм юп юр яг як ян яр яс яш

Found in the middle and at the end:

ай аы аю бм бс вм вр вс вт ві вю гм гі фу дт ей ец ею ея зд зл зм зь иж ий иф ич ию ия
йд ик йл йм йн йп йр йс йт йу йф йы йя кк кт лд лк лп лс лт лі ля мм мп мс мт нг нд нж нз
нк нн нр нс нт нч нш ня ёй ёп пп рб рв рг рд рк рқ рл рм рн рп рс рт ру рф рх рц рш рщ рі ру
рю тл тм тс тт тч уб уд уу уф ух фм фф цо шь щы ый ык ып ік ік іп іт іу ыб ыд ыз ык ым
ын ыс ыт ыф ыя юд юз юй юс ют юф яж яз яй як ял ям ят яу

Found at the beginning and at the end:

аш ГИ ПЭ

Found at the beginning and in the middle:

ағ ао аә әб әд әз әй әл әс әү бл бн бө бұ бу бю вл гв гл гн гр гү ғұ дә дв дұ еб
еҳ жә жү жә иа иб ии иі кб кв кн кө кү кы қо қө құ қү лә лұ лұ лю мә мө мұ мұ мә мю
нұ нұ оғ ох оц оя өб өг өж өм өр өт өш пл пұ пь рұ сә св сл сн сө сп сұ су сғ сх сц сю тә
тв тө тү тю уә ув уқ уо уч үғ үд үж үт үш үб үг үғ үд үж үм үп фи фл фу хе хи хл хо хр
ху хә це ци цу че чу шә шк шл шм шн шо шә шп шр шұ ығ ыд ыж іб әб әд әй әк әм әп әс әғ
юб юл юн яп ях

Found only at the end:
йю рч уя ұы ңы ще юы юи

Found only in the middle:

Found only at the beginning:

бы ву вы въ гѣ ꙗу дз дѳ дх дъ дѣ дю жѣ зѣ зꙗ кѣ кѣ кѹ кѧ мр нѣ ѡд пѣ пъ рѣ съ съ
уй ѹк хѹ цх шв шу шх ѩг эг эж эо эу эх эш юг

It is very important to know item distributions of vowel sounds in a word as in Kazakh language the accent falls on a final syllable. It considerably simplifies and reduces diphone base. For example, the sound **O** meets only at the beginning of a word and all diphones with sound participation **O** have no shock analogs.

Vowel sounds meeting at the beginning and in the middle of a word sound almost equally that also does diphone base of more compact.

Conclusion: Further development of the program involves automating the generation of other phonetic units (allophones, Trifonov, syllables). The criterion can serve not only phonetic (segment) text properties, but also communicative and syntactic component (for punctuation): for example, you can specify the proportion of interrogative sentences, or sentences containing nonfinal syntagma (by commas). Obviously, the efficiency of selection depends on the parameters of the reference body text: the more complete and varied it is, the more informative and compact it will be obtained on the basis of textual material.

It is also planned to automate the analysis of statistics of occurrence in the text of the speech units of different levels: phonemes, allophones, syllables, sound sequences of words, the development of automatic phonetic transcriptor Kazakh speech text to speech synthesizer Kazakh.

References

1. Shelepow V.U. Leksii o raspoznavanii rechi. – Donetsk: IPSHI «Nauka i osvita», 2009. – 196 p.

УДК 81'33

МНОГОФУНКЦИОНАЛЬНАЯ МОДЕЛЬ ТЮРКСКОЙ МОРФЕМЫ: БАЗОВЫЕ ФОРМАЛИЗМЫ

Сулейманов Джаведет Шевкетович д.т.н., директор, Институт прикладной семиотики Академии наук Республики Татарстан, Казанский федеральный университет.

E-mail: dvdt.slt@gmail.com

Гатшатуллин Айрат Рафизович к.т.н., зав.отделом, Институт прикладной семиотики Академии наук Республики Татарстан. E-mail: agat1972@mail.ru

Статья является продолжением описания многофункциональной модели тюркской морфемы, которая включает структурно-функциональную модель морфем, базу данных