

АВТОМАТИЧЕСКОЕ ЗАПОЛНЕНИЕ БД ПОРТАЛА ТЮРКСКОЙ МОРФЕМЫ С ПОМОЩЬЮ ПРОГРАММНОЙ ОБРАБОТКИ ДВУЯЗЫЧНЫХ СЛОВАРЕЙ

Аюпов М.М.

*Институт прикладной семиотики
Академии наук Республики Татарстан*

19 октября 2020 г.

При создании больших многоязычных баз данных важная роль принадлежит заполнению этих баз данных, так как готовых решений не существует и подготовка нужных данных занимает много времени.

7 языков для портала тюркской морфемы:

- татарский,
- киргизский,
- узбекский,
- крымско-татарский,
- казахский,
- чувашский,
- башкирский.

Russian	Semant
бегун	1192

Tatar	Semant
йөгерешче	1192

Uzbek	Semant
yuguruvchi	1192

Tatar	Uzbek
йөгерешче	yuguruvchi

Tatar	Kirgiz	Uzbek	Cirim tatar	Rus	Kazah	Bashkir
абайлылык, ихтираз,	аярдык, кыраакылык,	tuyg'unlik	ачыккозьлик, ихтият,	осторожность	сезімталдық	
абайлы, сагаюлы, сак,	абайлагыч, кыраакы		ачыккозь, мукъайт,	осторожный	қырағы	нак
йөгер, чап	жүгүр	chop, yugur	чап, ювур	бегать	шап	йүгерергә
бегемот, гиппопотам					сусиыр	
качкын	качкын, качуучу	qochkin	къачакъ, къачкын	беглец	қашқын	
жәһәтлек, йөгереклек,		ildamlik, tezlik	суръат, тезлик	быстрота	шапшаңдық	тизлек
качкан	качки			беглый	қашқын	
алгасак, алгыр, алгысак,	амалдуу, букта,	chopkir, ildam,	суръатлы, тез, чабик,	проворный	шапшаң	етез, өлгөр, сос, уңған
бегония			бал-къаймакъ	бегония	бегония	
йөгерешче, йөгерүче		yugurdak,		бегун	жүгіріш	
афәт, бәла, бәла-каза,	алат, апат,	balo, baxtsizlik,	бахытсызлыкъ, бея,	беда	сорлылык	бәлә
болама, буталчылык,	будуңчаң, иретсиздик,	sistemasizlik	къарышыкълыкъ,	беспорядок	жүйесіздік	бола
банкротлан, бөл,	жакырдан, жардылан,		сюмеле, факъырлаш,	бедность	кедейлен	бөлөргә, хәлһезләнергә
йолкышлык, мәхрүмлек,	бакырдык, бакырлык,	faqirlik, muhtojlik	ёкъсуллыкъ, фукъарелик	бедность	мұқтаждық	
фөкыйрь-фөкара, ярлы-				беднота	кедейлер	
барлыксыз, малсыз,	жакыр, жарды, кедей,	faqir	джарлы, ёкъсул, зюгюрт,	бедный	тіленшек	хәлһез, ярлы
бәхетсез, кайгы-	бактысыз, бакытсыз,		бедбахт, гъарип,	невезучий	сорлы	
бахыр, бичара,	бечара		байгъуш, бичаре	бедняга	бишара	меңкен
гидай	бакыр, кембагал			бедняк		фәкир
бот, сан	бут, сан	son		бедро	сан	бот

Язык	Количество слов
татарский	33066
казахский	18757
башкирский	2977
крымско-татарский	7065
киргизский	9686
узбекский	5433

Наиболее полные двуязычные словари в основном были созданы и изданы в прошлом веке. В связи с этим у них отсутствует электронная версия.

1. Найти отсканированную версию нужного словаря или, если нет версии с хорошим качеством, отсканировать.

2. Двужычны словарь распознаецца с
памошчу спецыяльных праграмм
распознавання.

АККУРА`ТНО нареч. ўкыпты,мўкыят,жынақы; ~ отвечать на письма хатка мўкыят жауап беру; ~
переписатьмўкыят көшіру; 2. разг. абайлап,байқап,жайлап,білдірмей,акырын; ~неситеабайлап
апарындар; ~узнайжайлап біл; 3. разг. жүйелі,үнемі,үздіксіз; ~ навешать больногосыратка үнемі
барып тұру.

3. С помощью программной обработки удаляется ненужная в дальнейшем информация.

АККУРАТНО ұқыпты, мұқият, жинақы 2. абайлап, байқап, жайлап, білдірмей, ақырын 3. жүйелі, үнемі, үздіксіз.

4. Обработанный словарь готовится для загрузки в базу данных.

АККУРАТНО	ұқыпты	1
АККУРАТНО	мұқият	1
АККУРАТНО	жинақы	1
АККУРАТНО	абайлап	2
АККУРАТНО	байқап	2
АККУРАТНО	жайлап	2
АККУРАТНО	білдірмей	2
АККУРАТНО	ақырын	2
АККУРАТНО	жүйелі	3
АККУРАТНО	үнемі	3
АККУРАТНО	үздіксіз	3

Казахский язык:

Экспериментальный словарь содержал **18757** строк.

Отсканированный pdf файл русско-казахского словаря содержал около 56000 словарных статей.

Готовый к загрузке в БД словник содержит более **117000** строк.

Возникают орфографические ошибки, если встречаются неуверенно распознанные слова и слова отсутствующие в словаре программы распознавания.

Отсутствие пробелов – распространенный случай при работе с распознанным текстом.

АБСОРБИРОВАТЬСЯ сов., несов. сорылу, сіңу, жұтылу.

АБСОРБЦИЯ ж. физ., хим. абсорбция; сіңірілу, жұтылу, сорылу.

АБСТРАГИРОВАНИЕ с. абстракциялаушылық.

АБСТРАГИРОВАТЬ сов., несов. что абстракциялау, дерексіздендіру, жалпылау.

АБСТРАГИРОВАТЬСЯ сов., несов. абстракциялану, дерексіздену.

АБСТРАКТНОСТЬ ж. дерексіздік, абстрактылық.

АБСТРАКТН||ый, -ая, -ое абстрактылы, дерексіз; ~ые понятия дерексіз ұғымдар; ~ое тождество абстрактылы тепе-теңдік.

Во время обработки словарей знак переноса строки обычно означает начало новой статьи. Но так же встречаются лишние знаки переноса строки.

ПЕРЕОБУЧА`ТЬ несов. см. переобучить. ¶

ПЕРЕОБУЧА`ТЬСЯ несов. 1. см. переобучиться; ¶

2. страд. от переобучать. ¶

ПЕРЕОБУЧЕ`НИЕ с. см. переобучить, переобучиться. - ¶

ПЕРЕОБУЧИ`ТЬ сов. кого-что қайта оқыту, жаңадан үйретіп шығару. ¶

В словарях часто встречаются толкования-отсылки, которые требуют дополнительной обработки.

ПЕРЕРАБО`ТА||ТЬСЯсов. 1. бойғасіну, денегетаралу; пища~ласьтамақбойғасінді; 2. · см.переработать5.¶

ПЕРЕРАБО`ТК||Аж. 1. см. переработать1–4; 2. разг. артыкістелгенуақыт; · уплатитьза~уартыкістелгенуақыгүшінтөлем; 3. (то, чтопереработано)істепшығарылғанөнім, · ұқсатылғанбұйым.¶

ПЕРЕРАБО`ТОЧН||ЫЙ·-ая, -оеұқсататын, өндейтін, қайтажасайтын; ~ыйпунктөндеу пункті.¶

ПЕРЕРАСПРЕДЕЛЕ`НИЕс. см.перераспределить.¶

Выводы:

- получилось сэкономить время и трудозатраты при подготовке данных для загрузки в базу данных портала тюркской морфемы;
- хорошо бы разработать универсальное приложение для обработки двуязычных словарей с возможностью настройки работы программы через интерфейс.

Спасибо за внимание!

Игътибарыгыз өчен рәхмәт!