

таггер шугнанских глаголов

Юрий Макаров, Макс Меленченко
первокурсники ОП ФиКЛ Школы лингвистики

Что мы делали

Преодолевали трудности:

- отсутствие единой орфографии: нужен был орфографический конвертер
- словарь Карамшоева не очень качественно оцифрован, нужно было приводить его в порядок
- разнообразие глагольных основ
- нерегулярность образования форм
- и прочие морфологические нюансы

Masc	Fem	Plur
andožd	andūvž	andūvž
indožd	indöžč	indöžč

примеры основ перфекта

А зачем всё это?

- есть проект по изучению шугнанского языка; в его рамках разрабатывается **корпус**, которому нужен **парсер**
- таггер — первый шаг на пути к парсеру

Словарь

п̄ин ч. дремать; прикорнуть, поспать немного; б. ии нахчир-
п̄ин к̄инум вздремну-ка (я) немножечко.

НАХЧӢР-ЇЕУ м. охота на кииков; б. wáz-ат Шерғозӣ-йат
Хиѳ-ām нахчӢр-ѳеу raw̄ун вад я, Шергази и Хиѳ собрались
охотиться на горных коз.

НАӢТИ: НАӢТУӢД, ж. наӢтойд; ед. наст. I-е л. наӢтийум,
сокр. наӢт̄им; 2-е л. мн. шб. наӢтийет, сокр. ш. тж. наӢтет;
3-е л. наӢт̄изд; перф. м. наӢтуӢч, ж. наӢтиц, мн. ч. наӢ-
тойч; инф. наӢтидов I) выходить (в разн. знач. извне, отку-
да-л., из чего-л.); ш. к̄ал̄онсух̄ӯб тар ваѳ ма̄-наӢти не вы-
ходи наружу с непокрытой головой; ш. ас х̄ац-анд-ен wāб чи-

Словарь: что мы с ним делали

- вручную вычитывали ошибки распознавания
- с помощью RegEx меняли систему глосс
- писали скрипты, извлекающие нужные для глагольной БД формы из распознанных словарных статей
- добавляли лексические значения к части глаголов

пйи ч. дремать; прикорнуть, поспать немного; б. йи нахчир-пйи кйи́нум вздремну-ка (я) немножечко.

НАХЧЫР-ЇЕУ м. охота на кииков; б. wáz-ат Шерғозы́-йат Хыф-а́м нахчыр-ѳеу rawу́н вад я, Шергази и Хыф собрались охотиться на горных коз.

НАХТИ:НАХТУЙД, ж. нахтойд; ед. наст. I-е л. нахтийум, сокр. нахтйм; 2-е л. мн. шб. нахтийет, сокр. ш. нахтет; 3-е л. нахтйэд; перф. м. нахтуйч, ж. нахтиц, мн. ч. нах-тойч; инф. нахтйдов I) выходить (в разн. знач. извне, отку-да-л., из чего-л.); ш. калбисуху́б тар ваç ма́-нахти не вы-ходи наружу с непокрытой головой; ш. ас ха́ц-анд-ен wáб чи-

Словарь: скрипт

```
def extract(gloss,place,source,sample):  
    if gloss in source:  
        n = source.index(gloss)  
        if n == -1:  
            pass  
        else:  
            sample[place] = source[n+1]  
    return sample
```

1. Находит помету (*gloss*) в строке (*source*) вида:
Praes ФИНУХС/ФИНУХЦ;
Past ФИНУХСТ/ФИНУХЦТ;
Perf финухсч/финухцч;
Inf бинухстow/финухцtow
медленно делать что-л.
2. Вычисляет её положение и добавляет её на место (*place*) в ячейке словаря словоформу после пометы.
3. `sample = ['0','1','2','3','4','5','6','7','8','9']`

Таблица глагольных форм и основ

PRS.M	PRS.F	PRS.3SG	PST.M	PST.FPL	PRF.M	PRF.F	PRF.PL	INF.M	INF.F	lemma
0	1	2	3	4	5	6	7	8	9	10
θāw	—	θöd	—	—	θuvž θūðž	θic	θavž θoðž	θid	—	—
naxti	—	—	naχtūjd	—	naχtūjž	naχtīc	—	naχtīd	—	<i>make cry</i>
pirmir parmir	—	pirmirt parmirt	pirmūd parmūd	pirmod parmod	pirmūyž parmūyž	pirmīyž parmīyž	pirmoyž parmayž	pirmird pirmīd pirmirt parmīd	—	—

Таблица глагольных форм и основ

101	wuqsēn	0	0	wuqsēnt/waqsēn/aqsēn	0	wuqsēnč/waqsēnč/aqsēnč	0	0	wuqsēnt/waqsēnt/aqsēnt	0	vomit	
102	wōj/woj	0	0	wōjt/wojt	0	wōjč/wojč	0	0	wōjt/wojt	0	cry_loudly	
103	wug	0	0	wugd	0	wōgu	0	0	wōgd	0	—	
104	gārō	0	0	gārōd	gāxt	0	gāxč	0	0	gāxt	0	—
105	gilneȳd	0	0	galneȳd	0	gilneȳʒ	0	0	gilneȳd/ʒalneȳd	0	—	
106	girz/čirz	0	0	girzd/čirzd	0	girzč/čirzč	0	0	1	0	grieve/mourn	
107	gīr	0	0	gīrt	0	gīrč	0	0	gīrt	0	agree	
108	gumbun	0	0	gumbunt	0	gumbōnč	0	0	gumbōnt	0	—	
109	gumor/gimor	0	0	gumort/gimort	0	gumorč/gimorč	0	0	gumort/gimort	0	authorize/delegate	
110	dawun	0	0	dawunt	0	dawōnč	0	0	dawōnt	0	—	
111	dak	0	0	dakt	0	dakč	0	0	dakt	0	dry_(about_throat)	
112	darjov	0	0	darjovd/darjevȳd	darjovȳd/darjevȳd	0	darjovȳʒ/darjevȳʒ	0	0	darjovȳd/darjevȳd	0	—
113	dām	0	0	dāmt	0	dāmč	0	0	dāmt	0	rotate	
114	deō	0	0	deōd	ded	0	deōʒ	0	0	ded	0	fan_the_fire
115	di/dē	0	0	dit/det	0	dīč/š/tž/dēč	0	0	1	0	enter	
116	diven	0	0	divent	0	divenč	0	0	divent	0	beat	
117	dives	0	0	divest	divixt	0	divixč	0	0	divixt	0	blow
118	diwen/dawen	0	0	diwent/dawent	0	diwenč/dawenč	0	0	1	0	show up/appear	

Код: глобальная структура



Код: конвертер орфографии

1.транскрипция из словаря (Д. Карамшоев)	Ūū	Ųų	Фф	Хх	Хх̣	Цц	Zz
2.иранистическая транскрипция (Д. Карамшоев)	Ūū	Ųų	Ff	Хх	Хх̣	Cc	Zz
3.транскрипция Плунгяна (см. 1, кроме i, ū, ě)	Ūū	Ěě	Фф	Хх	Хх̣	Цц	Zz
4.Зарубин (см.2)	Ūū	Ųų	Ff	Хх	Хх̣	Cc	Zz
5.Бахтибеков (см.1)	Ūū	Ųų	Фф	Хх	Хх̣	Цц	Zz
6.Umass (см.2, кроме много чего)	Uu	Oo	Ff	Хх	Хх̣	Tsts	Dzdz
7.Расторгуева (Эдельман) (см.2)	Ūū	Ųų	Ff	Хх	Хх̣	Cc	Zz
8.Morgenstierne (см.2, кроме Ёё и Óó)	Ūū	Ųų	Ff	Хх	Хх̣	Cc	Zz
9.Сердюченко Пахалина (см.2, кроме Ee)	Ūū	Ųų	Ff	Хх	Хх̣	Cc	Zz
10.Windfuhr (см.2)	Ūū	Ųų	Ff	Хх	Хх̣	Cc	
	29	30	31	32	33	34	35
IPA	u:	ũ	f	χ	x?	ts	dz
наша	ū	ö	f	χ	x	c	з
как в слове...	mūn 'ябл	göl 'немо	fawt 'сме	χūb 'хорс	жīn 'сини	сет 'глаз	awz 'ба

Код: конвертер орфографии

Naḡti-yēn

ortho.txt

Naḡti-jen

```
192
193 def orthoconv(text):
194     #конвертируем орфографию
195
196     with open('ortho.txt', 'r', encoding='utf-8') as file:
197         ortho = file.readlines()
198     goodlist = []
199     for line in ortho:
200
201         #с решёточки начинаются служебные строки в файле ortho
202         if not line.startswith('#'):
203             bad, good = line.split(' ')
204             while bad in text:
205
206                 #все плохие символы заменяются на хорошие
207                 good = good[0:len(good)-1]
208                 text = text.replace(bad, good)
209                 #print(bad+' > '+good)
210
211     return text
212
```

```
0 Tḡ C
1 tḡ č
2 Dz ž
3 dz ž
4 Sh š
5 SH š
6 sh š
7 #noncompounds
8 Æ Ā
9 æ ā
0 B B
1 b b
2 B V
3 v v
4 Γ G
5 Γ g
6 F B
7 F B
8 Y B
9 Y B
0 Ġ Y
1 í Y
2 Γ Y
3 Y Y
4 П D
```

Код: распознавание форм



```

589
590 def formdefinition(word, stem, y):
591     #это коммутатор, который выявляет, какая из основ найдена в слове, и перенаправляет программу к нужной функции (для verbfind)
592     #на вход принимается слово (токен текста) word, предполагаемая основа stem и индекс этой основы в матрице словаря y
593
594     #attributes – список глоссирований, который является результатом работы formdefinition
595     #здесь нужен именно список, потому что найденных глоссирований может быть больше одного
596     attributes = []
597
598     #если переданный индекс y равен индексу одного из стемов, идём к этому стему
599     #сначала проверяем на обычные формы
600     if word.startswith(stem):
601         if y == praes masc or y == praes femn:
602             attributes.append(isitpraestem(word, stem, y))
603         elif y == praes3sg:
604             attributes.append(isitpraes3sg(word, stem))
605         elif y == past masc or y == past fepl:
606             attributes.append(isitpasttnse(word, stem, y))
607         elif y == perf masc or y == perf femn or y == perf plur:
608             attributes.append(isitperftnse(word, stem, y))
609         elif y == inf masc or y == inf femn:
610             attributes.append(isitinfinite(word, stem, y))
611
612     #проверяем на стяжённые формы
613     if y == contracted:
614         attributes.append(isitcontract(word, stem))
615
616     #проверяем на отрицательные и условные формы
617     #если слово начинается на -ma- или -na-, а при удалении этих приставок найденная основа остаётся внутри слова
618     #то вызываем formdefinition ещё раз, но уже для слова без приставки (ДА!!! РЕКУРСИЯ!!! КРУТО ПРАВДА???)
619     #например, nahtijum не пройдёт, потому что stem nahti не сохраняется, если убрать -na-
620     #а вот naviđj пройдёт, потому что без -na- perfettoный stem viđj сохраняется, и formdefinition вызовется уже для viđj, чтобы
621     #определить его форму
622     if word.startswith('ma') or word.startswith('na'):
623         #если без приставки stem не сохраняется, мы предполагаем, что это стяжённая форма, и присваиваем индекс contracted
624         if not word[2:].startswith(stem):
625             if word.endswith('m') or word.endswith('et') or word.endswith('en'):
626                 y = contracted

```

```

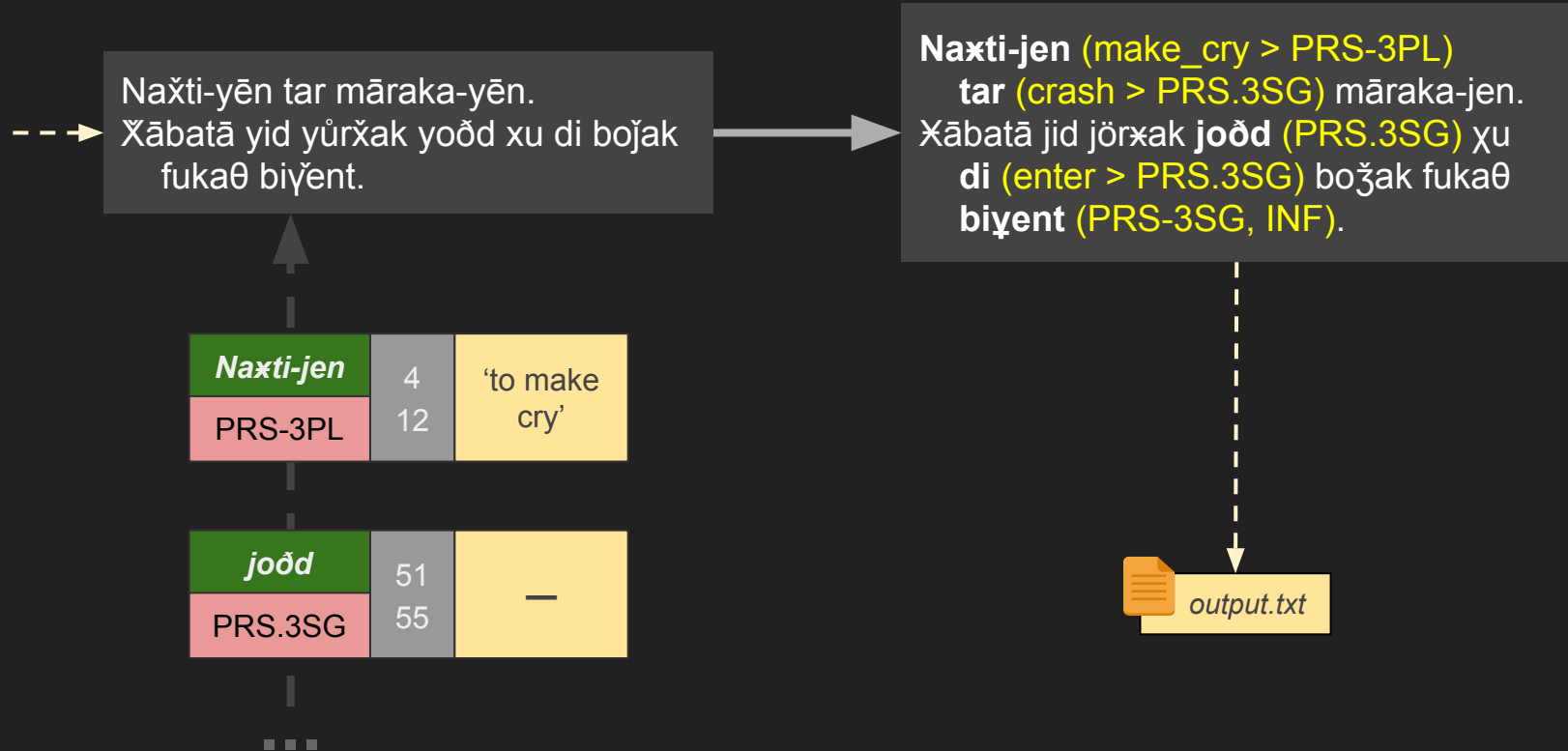
374
375 def isitpraestem(word, stem, y):
376     #функция проверяет, не является ли слово глагольной формой, образованной от praestem (для formdefinition)
377
378     #определяем род стема
379     if y == praesmasc:
380         gender = 'M'
381     if y == praesfemn:
382         gender = 'F'
383
384     #по умолчанию значение глоссирования – False, что сигнализирует о том, что формам этого стема токен не соответствует
385     attribute = False
386
387     #добавим к концу стема -j-, если он оканчивается на гласный
388     stem = jification(stem)
389
390     #проверяем, не является ли слово формой 3 лица ед. ч. презенса
391     if word.endswith('d') or word.endswith('t'):
392         flexias = ('d', 't')
393         for flexia in flexias:
394             if word == stem+flexia:
395                 attribute = 'PRS.'+gender+'-'+'+3SG'
396
397     #проверяем, не является ли слово одной из других форм презенса
398     if attribute == False:
399         flexias = {'um': '1SG', 'i': '2SG', 'ām': '1PL', 'et': '2PL', 'en': '3PL'}
400         for flexia in flexias:
401             if word == stem+flexia:
402                 attribute = 'PRS.'+gender+'-'+'+flexias[flexia]
403
404     #проверяем, не является ли слово отглагольным существительным
405     if attribute == False:
406         if word == stem+'i3':
407             attribute = 'AGENT_NOUN'
408
409     #возвращаем глоссирование (False, если ничего не найдено)
410     return attribute
411
412

```



```
parser.py vocab.txt ortho.txt graphic.py
697 #очистим слово от мусора – пробелов и висящих в начале и конце знаках препинаний – и будем использовать его для определения формы
698 wordnew = wordclean(word)
699
700 #если то, что осталось после очистки – нормальное слово, а не пустая строка, то продолжаем
701 if realword(wordnew):
702
703     #glossfoundsinword – переменная, которая для каждой итерации цикла собирает найденные в токене глоссирования, а потом пакует
704     #их и передаёт в glossboxes
705     glossfoundsinword = []
706
707     #для каждой лексемы в словаре:
708     for x in range(len(vocab)):
709
710         #для каждого индекса стема в лексеме:
711         for y in range(len(vocab[x])):
712
713             #если этот стем недостаточный, не трогаем его
714             if not nedostatochny_stem(vocab[x][y]):
715
716                 #для каждой вариации стема:
717                 for z in range(len(vocab[x][y])):
718
719                     #если мы нашли наш стем в токене
720                     if vocab[x][y][z] in wordnew:
721
722                         #отправляем токен, стем и его индекс на поиск формы и получаем attributes – глоссирование
723                         attributes = formdefinition(wordnew, vocab[x][y][z], y)
724
725                         #если какие-то глоссирования нашлись, т.е. список не пустой:
726                         if not attributes == []:
727                             for attribute in attributes:
728
729                                 #если можно указать лемму, то указываем её
730                                 if vocab[x][lemma][0] != '-':
731                                     attribute = vocab[x][lemma][0]+' > '+attribute
732
733                                 #складываем все глоссирования в glossfoundsinword
734                                 glossfoundsinword.append(attribute)
```

Вывод



Что получилось

github.com/iurmak/shughni

