

Persian Idioms

Collection and Identification in Texts

Vahid Ostovar

Unitex workshop

2014

Introduction

Persian idioms with the structure class *C1* (Maurice Groos, 1996)

Persian :	وحدید	تاب	خورد
Transliteration:	Vahid	tab	khord-past
Translation:	Vahid	swing	eat
Gloss:	Vahid had fun		

وحدید, *Vahid* is the subject (*N0*)

تاب, *tab* (*swing*) is the direct object (*C1*)

خورد, *khord* (*eats*) is the verb (*V*)

Project aim

- The main goal of this paper was to prepare a database of Persian idioms that could be used in the computational processing of that language.

Objectives

- Selection of web sources for the collection of idioms
- Building a database of Persian idioms
- Building FST tools for corpus exploring
- Corpus parsing and idioms candidate extraction
- Evaluation

Methodology

· Selection of web sources

- Queries on the web
- Using Google and AVG * browsers
- From January 2013 to May 2013
- Systematically survey of results, the first top ten pages per query

* AVG: <http://www.avg.com/eu-en/secure-search>

These are the most relevant queries that have been used.

Q1)

[فرهنگ اصطلاحات عامیانه زبان فارسی]

Farhange estelahate amiane zabane Farsi

“Persian language expressions dictionary”

Q2)

[اصطلاحات عامیانه و رایج به زبان فارسی]

Estelahate amiane va rayej be zabane Fars

“common expressions in Persian language”

Q3)

[فرهنگ اصطلاحات زبان فارسی]

Farhange estelahate zabane Farsi

“expressions dictionary of Persian language”

Several hits were discarded as not relevant, namely:

- Sites about English idioms and other languages in Farsi.
- Sites that did not contain enough idiomatic expressions and were really poor to be considered as a source.
- Many of them just give some general information about Persian idiomatic expressions and advertise dictionaries of idioms that will be published in the future.
- Websites that contain old Persian expressions.
- Sites appearing more than once in each query.

We retained 7 websites as a valid source and we collected 131, *C1* type idioms from them and the other idioms that is 233 different idioms, collected from two dictionaries of Persian idioms.

Building a database of Persian idioms (class C1)

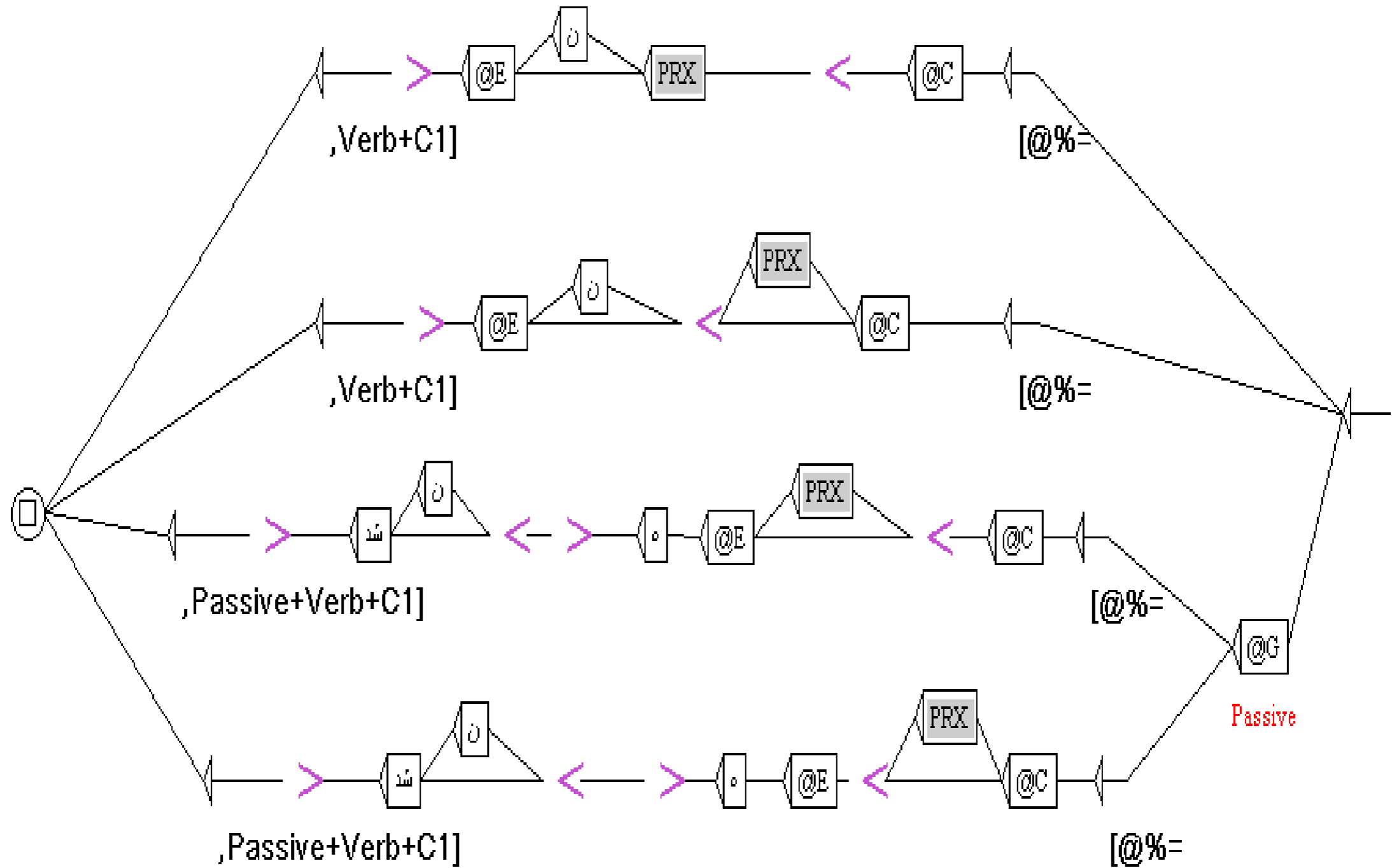
NO=:Hum	NO=:Nhum	C1	Det	Verb	Modify	Passive	Literal	Persian idioms	Transliteration	Gloss	Rough translation	Exact meaning	Reference
+	-	آب	<E>	آوردن	<E>	-	+	آب آوردن	ab avar dan	water bring	bringing water	being sick	http://www.scict.ir
+	-	آب	<E>	دادن	<E>	-	+	آب پس دادن	ab pas da dan	water give back	giving back the water	being generous -being unreliable	http://www.scict.ir
+	+	آب	<E>	رفتن	<E>	-	+	آب رفتن	ab raft an	water go	going water	became short- became low	http://www.scict.ir
+	-	آتش	<E>	کردن	<E>	-	+	آتش کردن	ata sh kar dan	fire do	making fire	shooting	http://www.scict.ir

Building FST tools for corpus exploration

Intersecting the tabular data with finite-state transducer (**FST**), using the Unitex* (3.0) corpus processing tool.

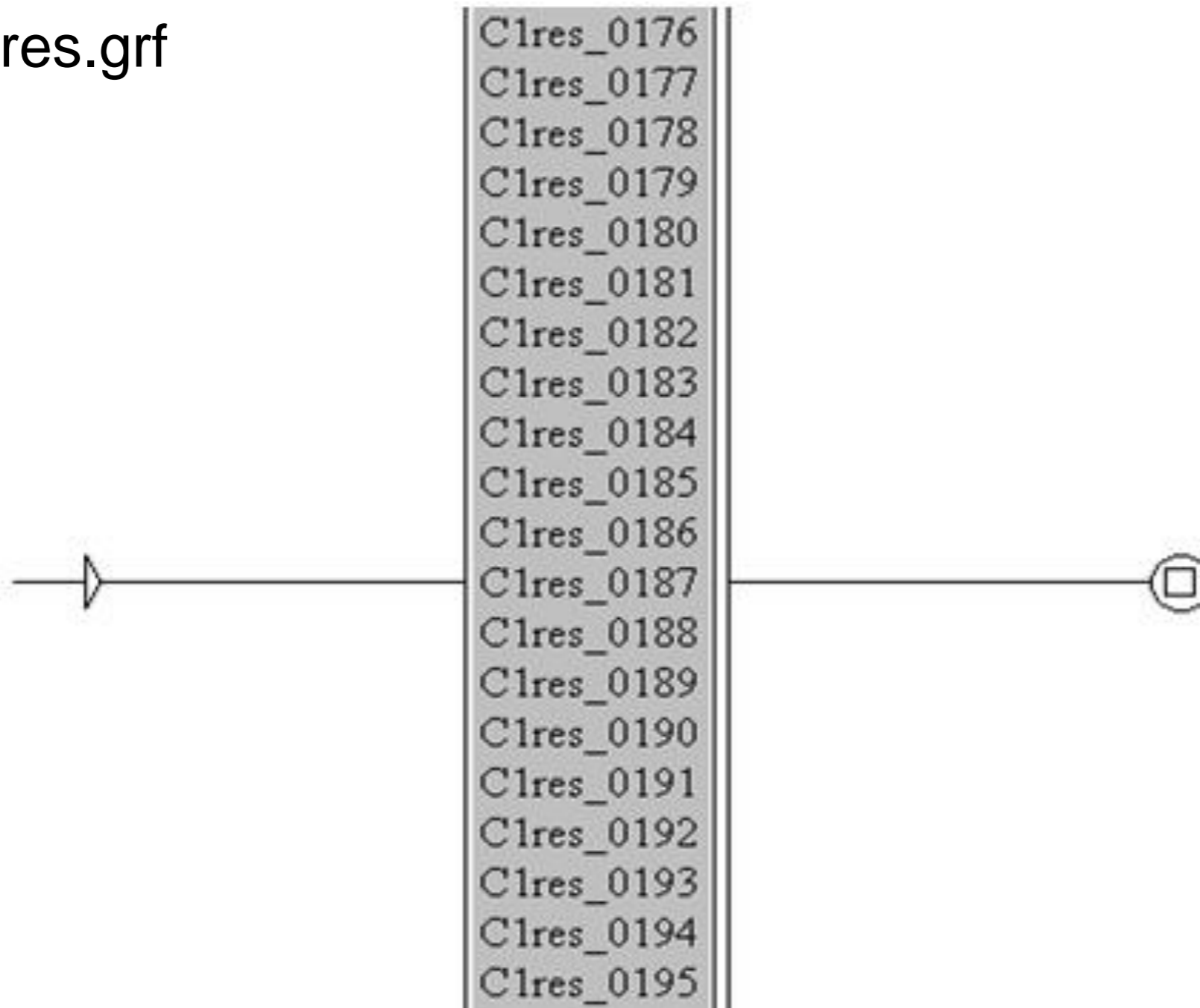
Unitex*: <http://www-igm.univ-mlv.fr/~unitex/index.php?page=0>

Reference graph



Intersecting the reference graph with the database

C1res.grf



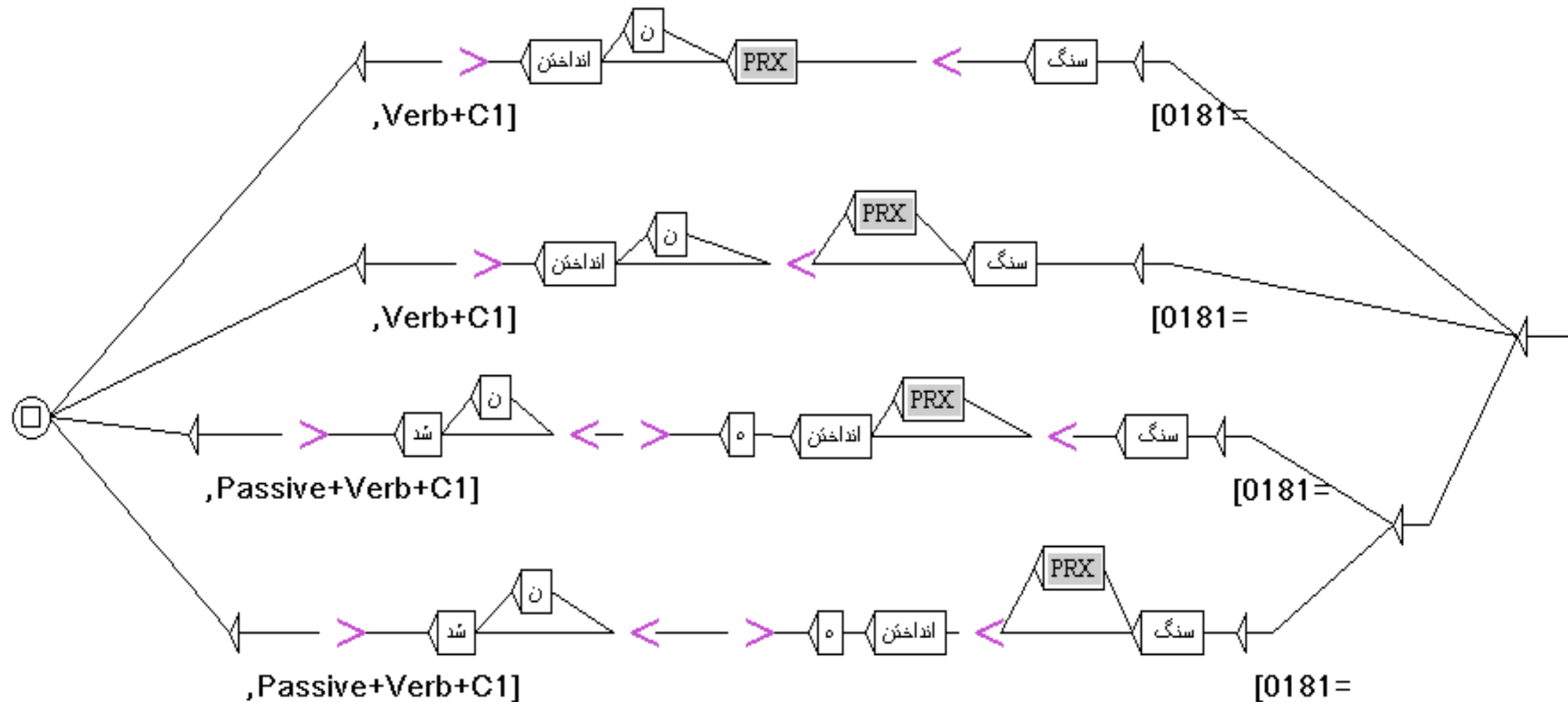
Persian : *وحد سنگ انداخت*

Transliteration: *Vahid sang andakht*

Translation: *Vahid is throwing stone*

Gloss: *Vahid is a trouble maker*

C1res.0181.grf



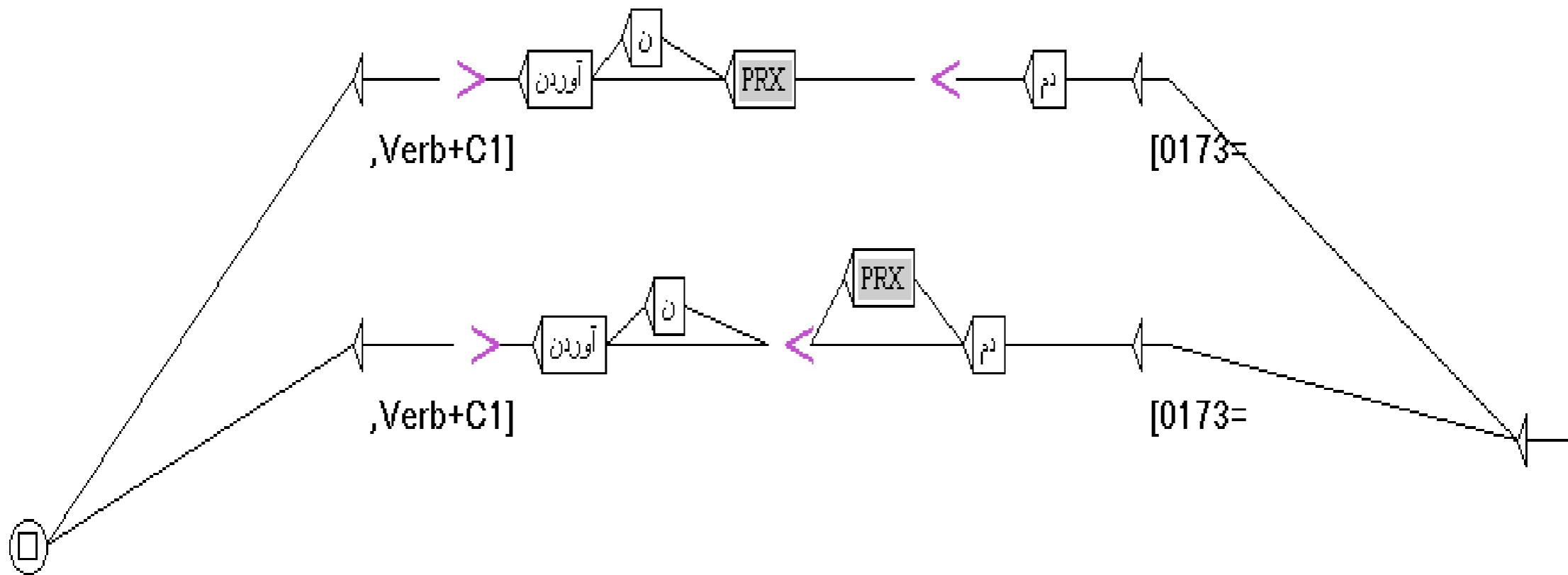
Persian: *وحدید دم در آورد*

Transliteration: *Vahid dom daravard*

Translation: *Vahid evolves tail*

Gloss: *Vahid becomes shameless*

C1res.0173.grf



Concordance from TEP corpus

م . 100701 اما مابل نبعتم که با سرو گوش [\[Verb+C1,دادن,آب=0002\]](#) زيادي براي همه مسئله درست کنم .
 اين به شما علامت میدم . 462879 اينجا با [\[Verb+C1,گرفتن,آتش=0005\]](#) دود غلیظ سياهي از خودتون ميدان
 شما چي فکر ميکنيد که اين اتفاق به وسيل [\[Verb+C1,گرفتن,آتش=0005\]](#) بک صفحه رخ داده بغه . 353933 م
 219165 من دوستانم را دیدم که جلوي چشمم [\[Verb+C1,گرفتن,آتش=0005\]](#) 219166 . کجا تو فابق . 219167
 د علفهاي زير فيروو ني که نو مزرعه بود ن [\[Verb+C1,گرفتن,آتش=0005\]](#) 127736 . اون هم سال 1935 بود .
 نه عزيزم تو هنوز هم ضابته اي . 322361 [\[Verb+C1,کردن,آب=0013\]](#) فقط بک امنياز اضافيه . 322362 ت
 بس کن هري سازمان کسايي را که عمه شون را [\[Verb+C1,کردن,آب=0013\]](#) نميفرسته آژکابان . 24428 از طرف
 . 63289 گمش نكي . 63290 من منظورشو از [\[Verb+C1,کردن,آب=0013\]](#) مارچ نفهميدم من فقط . 63291 کنت
 ار دوست داري اي . 99404 اينو بک بعير خوب [\[Verb+C1,آوردن,آب=0014\]](#) تا هميشه فروتنانه سفارش بده خي
 ريدم . 205391 ولي بين چه افتضاحي به [\[Verb+C1,آوردن,آب=0014\]](#) . 205392 چي از واگن ها دزد
 8 نه ببخشيد مامان من بک عقيداه محکم براي [\[Verb+C1,دادن,آب=0024\]](#) دارم . 148609 اگر نيتوني بري ت
 نگه هي حمله کنيم . 543588 جانگ بوگ آماده [\[Verb+C1,آوردن,آب=0025\]](#) خسارات بزرگه . 543589 اين شما
 درت فنوايي و نه قدرت رواني . 87287 براي [\[Verb+C1,آوردن,آب=0025\]](#) برابر قدرت با عظمت هداي واقعي
 46407 کمک رمي کمک . 46408 اميل شروع به [\[Verb+C1,آوردن,آب=0026\]](#) با لامپ . 46409 اميل به طرف من
 اهزنهائي که خودتون را به جايبرنس موکانگ [\[Verb+C1,آوردن,آب=0031\]](#) 424958 . همه جا هستند . 424959 بل
 اغش به بار من . 233216 ميدونستم که براي [\[Verb+C1,آوردن,آب=0031\]](#) قطعات باهوشيم . 233217 اين کاپربو
 اريس . 16022 دنبال من بيان لطفا . 16023 [\[Verb+C1,آوردن,آب=0032\]](#) اين دوتا هرم . 16024 خيلي نادر

Corpus parsing and idioms candidate extraction

The corpus TEP: Tehran English-Persian Parallel Corpus

- from Natural Language and Text Processing Laboratory of Tehran University
- 4 million tokens on each side
- Sentence Aligned

Processing with Unitex 3.0

556.234 sentence delimiters

15.166.987 (64.492 diff) tokens

4.485.147 (64.365) simple forms

3.239.250 (10) digits

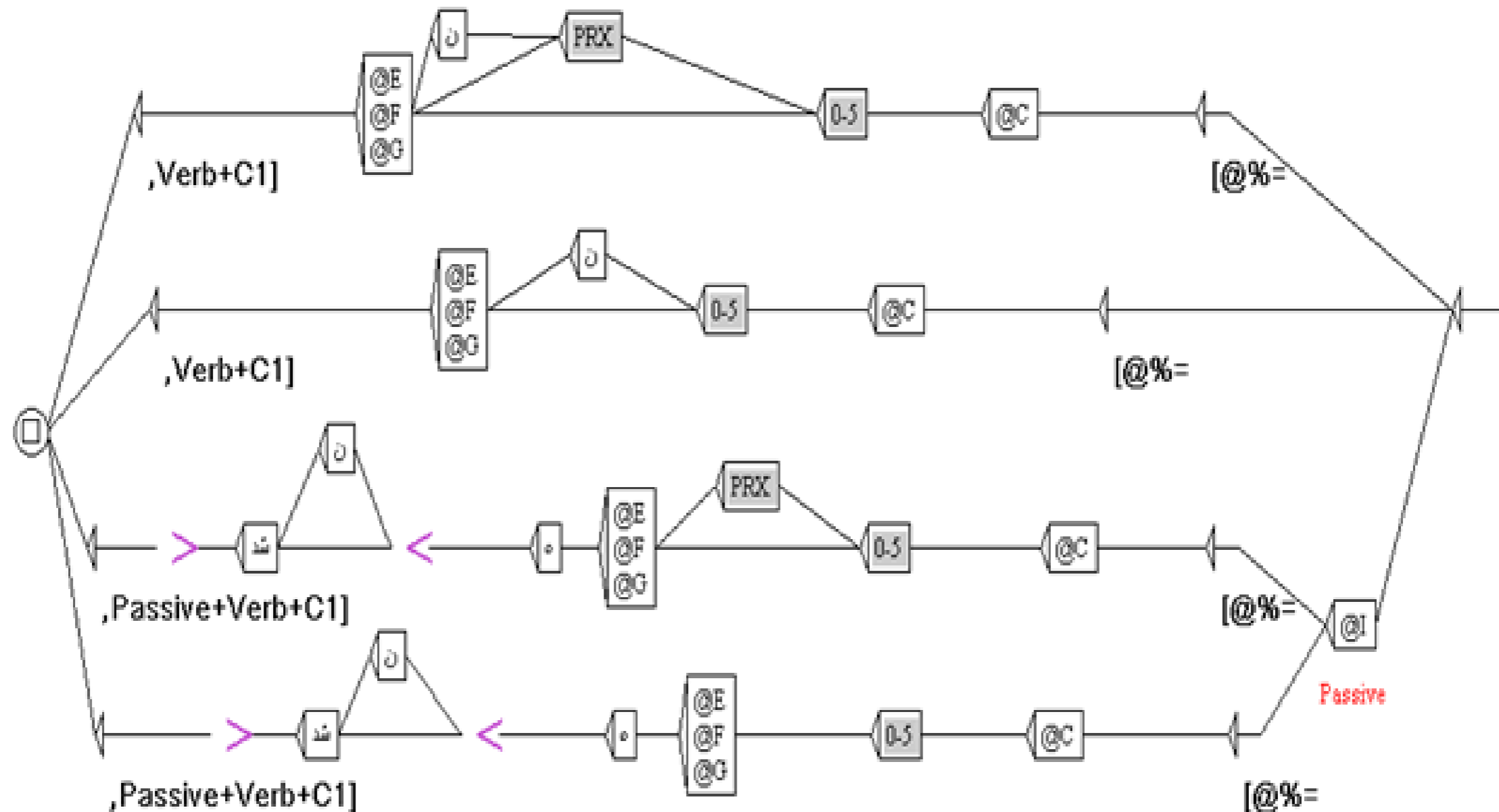
Evaluation

As the purpose of this project is to build a lexical database, the evaluation consist in three tasks:

- 1) Estimating the scope of the database
- 2) Determining the precision of the task of identifying the idioms, using the Unitex tools
- 3) Association measures

Estimation of idioms database scope

To estimate the scope of the database idioms, we produced this dummy FST to extract all the instances containing the verb and the lexical constant C1 from the text.



Database with the past and present stem of each idiom verb (the dummy FST is built according to this table)

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
NO=:Hum	NO=:Nhum	C1	Det	Verb	past tense	present ten	Modify	Passive	Literal	Persian idioms	Translitera tion	Gloss	Rough translation	Exact meaning	Reference
-	+	برق	<E>	افتادن	افتاد	افت	<E>	-	-	برق افتادن	bargh oftad	electricity fall	falling electricity	being too shiny and clean	http://www.scict.ir
+	-	افلاس	<E>	افتادن	افتاد	افت	<E>	-	-	افلاس افتادن	eflas oftad	Bankruptcy fall	falling bankruptcy	Getting caught in the poverty	http://www.sootak.ir
+	-	پس	<E>	افتادن	افتاد	افت	<E>	-	-	پس افتادن	pas oftad	back fall	falling back	fainting	Persian proverbs and idiomatic expressions book
+	-	تک	<E>	افتادن	افتاد	افت	<E>	-	-	تک افتادن	tak oftad	monad drop	dropping monad	being alone	Persian slang dictionary

Concordance of dummy FST

م . 100701 اما مایل نیستم که با سرو گوش [\[Verb+C1, دادن=0002\]](#) زیادی برای همه مسئله درست کنم .
 انواده ای که حکومت ازش میگذه . 7801 اونا [\[Verb+C1, دادن=0002\]](#) اونا ضبطونی نفونش میدن .
 تیش درست کنیم . 95581 می خواستم با بصرم [\[Verb+C1, کردن=0004\]](#) را تمرین کنیم . 95582 بافه
 هم کرد . 311825 آماده شو برای انفجار و [\[Verb+C1, روشن کردن=0004\]](#) همون جوری که گفتم . 311826
 این به شما علامت میدم . 462879 اینا با [\[Verb+C1, گرفتن=0005\]](#) دود غلیظ سباهی از خودتون میدن
 شما چی فکر میکنید که این اتفاق به وسیله [\[Verb+C1, گرفتن=0005\]](#) یک صفحه رخ داده بغه . 353933 م
 219165 من دوستانم را دیدم که جلوی جفام [\[Verb+C1, گرفتن=0005\]](#) 219166 کجا توفایق . 219167
 د علفهای زیر فیروونی که تو مزرعه بودن [\[Verb+C1, گرفتن=0005\]](#) 127736 اون هم سال 1935 بود .
 مغزت باید بیدار بغه . 87911 اونا من و [\[Verb+C1, جمع کردن=0011\]](#) 87912 کدوم کبر
 نه عزیزم تو هنوز هم ضایسته ای . 322361 [\[Verb+C1, کردن=0013\]](#) فقط یک امنیت از اضافه . 322362
 پس کن هری سازمان کسایی را که عمه شون را [\[Verb+C1, کردن=0013\]](#) میفرسنه آزکابان . 24428 از طرف
 63289 گمش نکنی . 63290 من منظورشو از [\[Verb+C1, کردن=0013\]](#) مارچ نفهمیدم من فقط . 63291 کنت
 ر مترقیه جادویی فورا اعزام شدند . 63215 [\[Verb+C1, کردن=0013\]](#) و حافظه اش را اصلاح
 نظر که بیام . 82599 خوبین بر شما سلامت [\[Verb+C1, کردن=0013\]](#) بهتر است . 82600 امروز میخو
 ار دوست داری ای . 99404 اینو یک پسر خوب [\[Verb+C1, آوردن=0014\]](#) تا همیشه فروتنانه سفارش بده چی
 ریدم . 205391 ولی ببین چه افتضاحی به [\[Verb+C1, آوردن=0014\]](#) 205392 چی از واگن ها دزد
 دنانون با ورود به این محاکمه . 8950 چند [\[Verb+C1, آوردن=0014\]](#) در این محاکمه آوردن . 8951 [\[Verb+C1\]](#)
 نگه می کنیم . 543588 جانگ بوگ آماده [\[Verb+C1, آوردن=0025\]](#) خسارات بزرگه . 543589 این فرما
 درت فنوایی و نه قدرت روانی . 87287 برای [\[Verb+C1, آوردن=0025\]](#) برابر قدرت با عظمت مدای واقعی

The results of the reference FSTs for idioms with the verbs کردن, *kardan* (to do) and زدن, *zadan* (to hit), with and without insertions.

Experiment	Number of entries (types)	Number of matches	Different idioms (types)
Entire database with insertion	364	1,754	115
Entire database without insertion	364	584	68
کردن, <i>kardan</i> (to do) with insertion	57	510	36
کردن, <i>kardan</i> (to do) without insertion	57	107	23
زدن, <i>zadan</i> (to hit) with insertion	47	84	12
زدن, <i>zadan</i> (to hit) without insertion	47	60	7

Concordances result of each dummy FSTs

	کردن, kardan (to do)			زدن, zadan (to hit)		
	TP	FP	S	TP	FP	S
FST w/ Ins	34	81	395	49	18	17
FST w/o Ins	53	54	0	42	18	0

Determine precision of identification by using the Unitex tools

	کردن, kardan (to do)	زدن, zadan (to hit)
FST w/ Ins	0.29	0.73
FST w/o Ins	0.50	0.70

Association measures

For the idioms that were captured by the FST using the entire lexicon-grammar of *kardan* (to do) and *zadan*, (to hit), we calculated two association measures; t-test and chi-square. (Manning and Schütze, 2005)

Idioms with the verb *زین*, *zadan* (to hit). T-test and Chi-square

Idiom	Example (transliteration) translation/gloss	T-test	Chi-square	t > p	$\chi^2 > p$
[0002]	جا زین (ja <i>zadan</i>) fearing to do something	2.949175	30.63058	1	1
[0004]	گند زین (gand <i>zadan</i>) sabotage	5.590053	2651.799	1	1
[0016]	قاطر زین (ghat <i>zadan</i>) becoming angry	1.402532	238.6465	0	1
[0023]	پرسه زین (parse <i>zadan</i>) walking without a goal	0.896753	7.804938	0	1
[0026]	نیش زین (nish <i>zadan</i>) Squibbing	1.970058	259.7837	0	1
[0035]	قاپ زین (ghap <i>zadan</i>) sudden stealing	1.405453	319.5216	0	1
[0045]	جوش زین (josh <i>zadan</i>) be concerned	1.288642	18.74106	0	1

Idioms with the verb کردن, *kardan* (to do)"T-test and Chi-square

Idiom	Example (transliteration) translation/gloss	T-test	Chi-square	t > p	$\chi^2 > p$
[0001]	اتش کردن (atash <i>kardan</i>) shooting	-9.36316	8.575242	0	1
[0003]	باد در کردن (bad <i>darkardan</i>) farting	2.629322	21.82092	1	1
[0005]	تابلو کردن (tablo <i>kardan</i>) revealing something	0.641734	1.164957	0	0
[0007]	توش کردن (toesh <i>kardan</i>) becoming angry	1.78106	29.36749	0	1
[0008]	جفت کردن (joft <i>kardan</i>) fearing	-0.96383	0.479417	0	0
[0010]	کف کردن (kaf <i>kardan</i>) getting excited	-0.43306	0.132633	0	0
[0012]	باد کردن (bad <i>kardan</i>) getting pregnant	2.629322	21.82092	1	1
[0015]	فاطی کردن (ghati <i>kardan</i>) becoming angry	2.951347	544.3165	1	1
[0016]	کلید کردن (kilid <i>kardan</i>) being bête noire	0.734618	2.060893	0	0
[0017]	تخم کردن (tokhm <i>kardan</i>) greatly fear	-1.5742	0.975672	0	0
[0018]	ردیف کردن (radif <i>kardan</i>) making everything alright	1.311004	23.86749	0	1
[0021]	پف کردن (pof <i>kardan</i>) feeling self-respect	2.357399	149.8065	0	1
[0022]	توش کردن (torsh <i>kardan</i>) huffing	1.78106	29.36749	0	1
[0027]	پنجر کردن (panchar <i>kardan</i>) being too tired	1.292239	19.62157	0	1
[0036]	اب باز کردن (ab baz <i>kardan</i>) happening something wrong	-5.20311	9.252444	0	1
[0038]	افاقه کردن (efaghe <i>kardan</i>) having good effect	0.986731	74.36311	0	1
[0040]	بچه درست کردن (bache dorost <i>kardan</i>) reproduction	-52.302	52.0676	0	1
[0041]	دراز کردن (darz <i>kardan</i>) make known something private	0.814233	3.616851	0	0

Comparison of T-test and Chi-square

	T & χ^2 (1 1)	\sim T & χ^2 (0 1)	T & $\sim \chi^2$ (1 0)	\sim (T & χ^2) (0 0)
کردن (kardan) to do	4	11	—	8
زدن (zadan) to hit	2	5	—	—

Conclusion

- Results indicate that depending on the verb, namely, if the verb has a more grammatical status or is a full verb, it may be necessary to adjust insertion windows between the key elements of the idiom.
- Precision varies depending on the verb.
- Association measures are valid strategies to capture these types of idioms; however, their precision depends on the idioms' verb, as the presence of support verbs may vary the results.

Future work

- Enhance the syntactic complexity of the reference graphs in order to capture other syntactic patterns such as *permutation*
- Plug the existing graphs with the Persian dictionary of inflected forms
- Extend the collection of idioms by using other sources, augmenting the lexical coverage of the lexicon-grammar of the idioms.
- Extend the collection of idioms to other types of Persian idioms

Thank
you