Чанкинг именной группы в свете разработки синтаксического подкорпуса Национального корпуса калмыцкого языка (предварительные замечания)*

Chunking of a Nominative Group in the Context of Development of the Syntactical Sub-Corpus of the Kalmyk National Corpus (preliminary notes)

В. В. Куканова (V. V. Kukanova)¹

¹ кандидат филологических наук, директор, Калмыцкий научный центр РАН (г. Элиста). E-mail: vika.kukanova@gmail.com.

Ph. D. in Philology (Candidate of Philological Sciences), Director, Kalmyk Scientific Center of the RAS (Elista). E-mail: vika.kukanova@gmail.com.

Аннотация. В данной работе рассматриваются проблемы чанкинга именной группы в свете разработки подкорпуса Национального корпуса калмыцкого языка. На сегодняшний день корпус снабжен морфологической и синтаксической разметками, другие виды аннотирования до настоящего времени не разрабатывались, в их числе синтаксический подкорпус. В статье описываются именные группы, составные элементы именной группы, а также алгоритм анализа, критерии, границы.

Ключевые слова: Национальный корпус калмыцкого языка, синтаксический подкорпус, чанкинг, именная группа.

Abstract. The work considers the problems of chunking of a nominative group in the context of development of a sub-corpus of the Kalmyk National Corpus. At present the Corpus is provided with morphological and syntactical tracking mechanisms, other types of data annotation — including a syntactical sub-corpus — have not been elaborated so far. The article describes nominative groups, component elements of a nominative group as well as the analysis algorithm, criteria and borders.

Keywords: Kalmyk National Corpus, syntactical sub-corpus, chunking, nominative group, word combination.

^{*} Статья подготовлена при финансовой поддержке РГНФ (проект «Структурно-вероятностный синтаксис монгольских языков», № 15-04-00553).

Национальный корпус калмыцкого языка является аннотированной системой текстов на калмыцком языке современного периода, преимущественно второй половины XX в. Его использование в качестве материала исследования и для составления лексикографических работ открывает широкие возможности для лингвиста, занимающегося изучением проблем в области лексики и грамматики. Однако на данный момент корпус снабжен только морфологической и синтаксической разметками, другие виды аннотирования до настоящего времени не разрабатывались [Куканова 2014; Куканова, Каджиев 2014]. Одной из актуальных задач является развитие корпуса и в направлении аннотирования синтаксической структуры калмыцкого предложения.

Существует два способа. Прежде всего следует указать на то, что в создании синтаксического корпуса с использованием нисходящего и восходящего анализа, существенное различие сводится к тому, как продвигается анализ: от «корня» к «листьям» (от главного к зависимому) или «листьев» к «корню» (от зависимого к главному). Эффективность применения того или иного подхода к калмыцкому предложению в настоящий момент не известна, поскольку отсутствуют какие-то практические данные, на которые можно было бы опереться в выборе подхода к анализу. Но для обоих подходов первоначальным этапом должна быть выработка синтаксических формальных правил, которые используются для проверки гипотез. Применение формального подхода, очевидно, также не решит проблему эффективности автоматического синтаксического разбора, поскольку при контекстно-свободной грамматике два варианта ус уух 'пить воду' и !!! ширә уух 'пить стол' будут правильными с точки зрения реализованной грамматической связи, при контекстно-зависимом подходе, т. е. с применением правил семантических ограничений, второй вариант будет признан неправильным с точки зрения семантики. Но и для разработки парсера, основанного на правилах контекстно-зависимой грамматики, необходимо описание формальных правил, на которые уже впоследствии наслаиваются правила семантики.

В данной работе мы сфокусируем внимание на синтаксически связанных непересекающихся фрагментах предложения — именной

группе. В калмыцком языке традиционно выделяются именные и глагольные группы, в их числе содержится большая группа фрагментов предложения, которые носят аналитический характер. В данной работе мы не рассматриваем их, поскольку было принято решение начать с описания «простых» именных сочетаний.

Принято выделять базовую именную группу, под которой понимается, что главный элемент — это имя, что от него зависят другие элементы или подчиняются ему и что данная конструкция носит элементарный характер, т. е. не содержит рекурсивный элемент. Другими словами, соблюдается правило нерекурсивности. Однако в калмыцком языке требование рекурсивности мы несколько изменяем, как, например, в случае конструкции NP -> Gen + N, где N=Gen рассматриваем как элемент, реализующий определительные отношения в первую очередь, а во вторую — как элемент, управляемый главным словом. Во многом в трактовках синтаксической связи мы опирались на классическую работу Ю. С. Маслова «Введение в языкознание» [1987].

Именная группа может состоять из вершины (обязательного компонента), премодификатора и постмодификатора (необязательных компонентов). Главный элемент, или вершина, выполняет функцию субстантива, в его качестве могут выступать:

- 1) существительное (N);
- 2) местоимение-существительное (Pron-N);
- 3) местоимение-прилагательное (Pron-Adj);
- 4) прилагательное (Adj);
- 5) числительное (Num);
- 6) причастие¹ (Ptcpl);
- 7) функционально близкая к именной группе конструкция (словосочетание), вершиной которого являются неопределенные количественные слова.

Для их описания мы использовали следующие условные символы:

- NP именная группа,
- c case,
- num number,

¹ Субстантированное причастие.

- Noun существительное,
- Pron местоимение,
- Adj прилагательное,
- Num числительное,
- Adv наречие,
- -> знак раскрытия,
- = знак совпадения в той или иной граммеме,
- — зависимый элемент, расположенный слева от вершины,
- - знак примыкания к левому элементу,
- — знак примыкания к правому элементу,
- х любой элемент граммемы,
- + элемент присоединения,
- Туре родовое понятие, тип, класс предметов,
- CommonNoun имя нарицательное,
- f в функции;
- Ass комитатив, совместный падеж;
- Ptcpl причастие;
- Name имя;
- Surname фамилия;
- Patronic отчество;
- N. сокращённый инициал имени лица;
- Post послелог;
- Conj союз;
- Тегт предельный падеж;
- Ins творительный падеж;
- Abl удалительный падеж;
- Dat дательный падеж.

Полужирным выделены вершины в группе.

Вершина может иметь различные виды зависимых, которые ее определяют. По-другому их называют модификаторами. В калмыцком языке они могут стоять только в препозиции, поэтому их называют премодификаторами, в качестве которых могут выступать:

- 1) прилагательные (байн күн 'богатый человек', хурдн мөрн 'быстрая лошадь' и т. д.);
- 2) местоимения-прилагательные (эн гер 'этот дом', тер хаша 'тот забор; та ограда');

- 3) числительные-прилагательные (*негдгч класс* 'первый класс', *хойрдгч машин* 'вторая машина');
- 4) именные адъюнкты генетивные, совместные и т. д. конструкции (ухата көвүн 'умный мальчик', хальмгин үкр 'корова калмыка');
- 5) одиночное причастие или причастный оборот, выступающий в роли релятивных оборота (Дөрвн сард хурһ көкүлсн хөөдиг 'овцематкам, кормившим в течение четырех месяцев ягнят');
- 6) послеложные группы (герин тускар 'о доме');
- 7) конструкции типа отглагольное имя + существительное (умилhна дегтр 'книга для чтения');
- 8) отрицательные конструкции (*көл уга ширә* 'стол без нож-ки').

Рассмотрим виды именных групп.

- 1. Именные. Здесь вершиной является часть речи в функции субстантива, которая может иметь одиночное или развернутое определение:
 - а) одиночное имя (вершина);
 - б) существительное в косвенном падеже + имя (иколын журнал 'школьный журнал', мөрт күн букв. 'человек с лошадью');
 - в) послеложная группа: имя + послелог (*дурна тускар* 'о любви');
 - г) отглагольное имя + имя (*саалhна аппарат* 'доильный аппарат');
 - д) имена собственные:
 - сочетание имен собственных;
 - несобственные наименования сочетание имени нарицательного и собственного.

Шаблон	Пример	
Именные		
NP -> Noun	ширә 'стол'	
NP -> Noun (c=Ass, num=x)	бүшмүдтэнь 'та, что в платье'	
NP -> Adj (c=n, num=x)	цецн 'мудрый'	

NP -> Num	хоюрн 'обе/оба'	
NP -> Ptcpl (c=Nom)	сансн 'думающий'	
NP -> Num Noun (c=x)	цөөкн яман 'несколько коз'	
NP -> Num - Adj - Noun (c=x)		
	KO3'	
$NP \rightarrow Noun (c=x) + Noun (c=x)$	күүкн сурһульч 'девочка-школьница'	
NP -> Noun (c=x) - Neg	шил уга букв. 'стекол нет'	
$NP \rightarrow Noun (c=Gen) \leftarrow Noun$	эцкин гер 'дом отца'	
(c=x, num=x)	эцкин герәс 'из дома отца'	
$NP \rightarrow Noun (n, c=Ass) \leftarrow Noun$	машитә залу букв. 'мужчина с машиной'	
(c=x, num=x)	машитә залуһур 'к мужчине с машиной'	
$ NP \rightarrow Noun (c=Dat, n=x) \leftarrow$	сурһульчнрт заавр 'указания учащимся'	
Noun (c=x, num=x)	медәтирт дөңг 'помощь старикам'	
	үүртән иџлһн 'надежда на друга'	
	увлд белдлин 'подготовка к зиме'	
NP -> Adj - Noun (c=x, num=x)	цаћан чирә 'светлое лицо'	
	цаћан чирәтә 'со светлым лицом'	
$NP \rightarrow Adj \mid Adj \mid Noun (c=x,$	сәәхн цаһан чирә 'красивое белое лицо'	
num=x)	сәәхн цаһан чирәтә 'с красивым белым	
	лицом'	
NP -> Ptcpl Noun (c=x, num=x)	ирсн күн 'пришедший человек'	
$NP \rightarrow Post \leftarrow Noun (c=x, num=x)$	ширә деер 'на столе'	
$ \text{NP} -> \text{Noun (c=Gen)} \leftarrow \text{Noun} $	услһна аппарат 'питьевой аппарат'	
(c=x, num=x)	услина аппарато 'к питьевому аппарату'	
$ \text{NP} -> \text{Noun } (c = \text{Acc}) \leftarrow \text{Noun} $	чидл үзүллһн 'демонстрация силы'	
(c=x, num=x)		
$NP \rightarrow Adv \leftarrow Noun (c=Ass)$	йир күчтә залу 'очень сильный мужчина'	
	а собственные	
NP-> Name Surname Patronic	Иван Сергеевич Тургенев	
(c=x, num=x)		
NP -> Adj Name Surname	алдр Лев Николаевич Толстой 'великий	
Patronic (c=x, num=x)	Лев Николаевич Толстой'	
NP -> Pron-Adj - Name мана Лев Николаевич Толстой 'наш Ле		
Surname Patronic (c=x, num=x) Николаевич Толстой'		
NP -> Pron-Adj Adj Name мана алдр Лев Николаевич Толстой 'на		
Surname Patronic (c=x, num=x)	великий Лев Николаевич Толстой'	
NP -> N.P. Surname (c=x,	И.С. Тургенев	
num=x)		

$NP \rightarrow Adj \mid N.P. Surname (c=x,$	алдр Л.Н. Толстой 'великий Л.Н. Толстой'	
num=x)		
$NP \rightarrow Pron-Adj \mid Adj \mid N.P.$?. мана алдр Л.Н. Толстой 'наш великий	
Surname (c=x, num=x)	Л.Н. Толстой'	
NP -> Pron-Adj N.P. Surname	мана Л.Н. Толстой 'наш Л.Н. Толстой'	
(c=x, num=x)		
NP -> Patronic (Name (c=Gen))	Манжен Нимгр 'Манджиев Нимгир'	
Name (c=x, num=x)		
NP -> Adj Patronic (Name	алдр Манжсин Нимгр 'великий Манджиев	
(c=Gen)) Name (c=x, num=x)	Нимгир'	
NP-> Pron-Adj Adj Patronic	мана алдр Манжин Нимгр 'наш великий	
(Name (c=Gen)) Name (c=x,	Манджиев Нимгир'	
num=x)		
	мана Манжин Нимгр 'наш Манджиев	
(Name (c=Gen)) Name (c=x,	Нимгир'	
num=x)		
NP -> N. Surname (c=x, num=x)		
$NP \rightarrow Adj + N. Surname (c=x,$	алдр Л. Толстой 'великий Л. Толстой'	
num=x)		
NP -> Pron-Adj + Adj + N.	мана алдр Л. Толстой 'наш великий	
Surname (c=x, num=x)	Л. Толстой'	
- •	мана Л. Толстой 'наш Л. Толстой'	
(c=x, num=x)		
NP -> Noun (c=Nom) + Name	күүкн Цаһан 'девочка Цагана'	
(c=Nom)		
NP-> Pron-Adj - Noun (c=Nom)	чини күүкн Цаһан 'твоя девочка Цагана'	
+ Name (c=Nom)		
NP -> Pron-Adj Adj Noun	<i>чини сәәхн күүкн Цаһан</i> 'твоя красивая	
(c=Nom) + Name (c=Nom)	девочка Цагана'	
NP -> Adj Noun (c=Nom) +	<i>сәәхн күүкн Цаһан</i> 'красивая девочка	
Name (c=Nom)	Цагана'	
$NP \rightarrow Name (c=x) + Noun (c=x)$	<i>Цаћан күүкн</i> 'девочка Цагана'	
	<i>Цаһанла күүкнлә</i> 'с девочкой Цаганой'	
$NP \rightarrow Adj$ - Name (c=x) + Noun	сәәхн Цаһан күүкн 'красивая девочка	
(c=x)	Цагана'	
NP -> Name (c=x, num=x)	Бадм 'Бадма'	
NP->Adj Name (c=x, num=x)	өндр Бадм 'высокий Бадма'	
NP -> Pron-Adj Name (c=x, тана Бадм 'Ваш Бадма'		
num=x)		

$NP \rightarrow Noun(n, c=Ass) \leftarrow Name$	машитэ Бадм 'Бадма с машиной'	
(c=x, num=x)		
NP -> Surname (c=x, num=x)	Нимгиров 'Нимгиров'	
$NP \rightarrow Noun (n, c=Ass) \leftarrow$	машитэ Нимгиров 'Нимгиров с	
Surname (c=x, num=x)	машиной'	
$NP \rightarrow Noun (n, c=Gen) \leftarrow Name$	Нинан Бадм 'Бадма Нины'	
(c=x, num=x)		
$NP \rightarrow Type (c=x, n=a) +$	нойон Тундутов 'нойон (князь) Тундутов'	
Surname (c=x, num=x)		
$ NP \rightarrow Adj Type (c=x, num=x)$	алдр нойон Тундутов 'великий нойон	
+ Surname (c=x, num=x)	(князь) Тундутов'	
NP -> Type (c=x, num=x) +	нойон Данзан Давидович Тундутов	
Name Patronic Surname (c=x,	'нойон-князь Данзан Давидович	
num=x)	Тундутов'	
$ NP \rightarrow Adj Type (c=x, num=x)+$	алдр нойон Данзан Давидович Тундутов	
Name Patronic Surname (c=x,	'великий нойон-князь Данзан Давидович	
num=x)	Тундутов'	
$ NP \rightarrow Type (c=x, n=x) + Name $	нойон Данзан Тундутов 'нойон-князь	
Surname (c=x, num=x)	Данзан Тундутов'	
$ NP \rightarrow Adj Type (c=x, num=x)$	алдр нойон Данзан Тундутов 'великий	
+ Name Surname (c=x, num=x)	нойон-князь Данзан Тундутов'	
NP -> Type (c=x, num=x) +	нойон Данзан 'нойон-князь Данзан'	
Name (c=x, num=x)		
$ NP \rightarrow Adj Type (c=x, num=x)$	алдр нойон Данзан 'великий нойон-князь	
+ Name (c=x, num=x)	Данзан'	
NP -> Type (c=Nom) +	<i>hол Иж,</i> л 'река Волга'	
CommonNoun (c=Nom)		
NP -> Adj Type (c=Nom) +	гүн нүр Байкал 'глубокое озеро Байкал'	
CommonNoun (c=Nom)		
NP -> CommonNoun (c=x) +	Иж,л һол 'Волга-река'	
Type (c=x)	<i>Ижд hолд</i> 'в / на Волге-реке'	
NP -> Adj CommonNoun	гүн Байкал нүр 'глубокое Байкал-озеро'	
(c=x) + Type (c=x)		

- 2. Адъективные. Здесь вершиной является прилагательное, которое понимается условно и включает в себя слова, которые могут вступать в определительные отношения с вершиной.
- а) имя, модифицирующее прилагательное, + прилагательное + имя (малар + байн + күн 'человек, богатый скотом');

- б) наречие, модифицирующее прилагательное, + прилагательное + имя ($u\ddot{u}p + c\partial h + \kappa \gamma \gamma \kappa h$ 'очень хорошая девочка');
- в) наречие, модифицирующее имя в совместном падеже, + имя в совместном падеже + имя ($\~uup + yxama + \kappa \theta в y h$ 'очень умный мальчик');
- Γ) имя в аблативе + прилагательное + имя (*сәәнәс сән сурһульч* 'лучший из лучших ученик').

Модификатор-адъектив может появляться в следующих конструкциях:

- а) одиночное прилагательное + имя (көк теңгр 'синее небо');
- б) однородные прилагательные + имя (байн + xap + ych 'густая черная кровь');
- в) причастие или причастный оборот ($hapch \ \theta \partial p$ 'день рождения', $uup \partial h \partial h \ mypk \partial \epsilon \ kpem$ 'крем для лица').

Шаблон	Пример	
Адъективные		
$NP -> Noun (c=Inst) \leftarrow Adj$	алтар байн күн 'золотом богатый	
	человек'	
$NP \rightarrow Noun (c=Term) \leftarrow Adj$	өвдгцэ гүн 'глубокий до колена'	
$NP \rightarrow Noun (c=Abl) \leftarrow Adj$	<i>эмтнлэ сэн</i> 'с людьми хороший'	
$NP \rightarrow Noun (c=Dat) \leftarrow Adj$	көвүнд күнд 'тяжелый [для] мальчика'	
$NP \rightarrow Adv \leftarrow Adj$	йир өндр бахн 'очень высокий столб'	
$NP \rightarrow Noun (c=Abl) \leftarrow Adj$	мууhас му күн 'худший из худших	
	человек'	
$NP \rightarrow Neg \leftarrow Adj$	хурдн уга 'небыстрый'	

Шаблон	Пример	
Местоименные		
NP -> Pron-N (c=n, num=x)	чи 'ты'	
NP -> Pron-Adj (c=n, num=x)	манахс 'наши'	
NP -> Adj Pron-Adj	күндтә манахс 'наши'	

4. Количественные. Модификатором вершины служит числительное или числа: числительное (или последовательность числительных) + имя (хойрдгч + көдлмш 'вторая работа', арвн + арслң 'десять рублей', арвн долан арслң 'семнадцать рублей').

Шаблон	Пример		
Количественные			
NP -> Num Noun (c=x)	хойр альмн 'два яблока' негдгч дивизь 'первая дивизия'		
NP -> Num Num Noun (c=x)	нәәмн зун жирн хойр арслң '862 рубля'		

5. Сочиненные — группа однородных именных групп, возможно объединённых союзом (например, *Бадм болн Цаћан* 'Бадма и Цагана'). Возможен случай с личным местоимением (например, *би болн Саша* 'я и Саша').

Шаблон	Пример	
Сочиненные		
$NP ext{->} extbf{Noun} (c=x) + Conj + extbf{Noun} $ $(c=x)$ $\kappa \Theta B \gamma H \ \delta O \pi H \ \kappa \gamma \gamma \kappa H \ \epsilon$ мальчик и де		
$NP \rightarrow Pron (c=x) + Conj + Pron (c=x)$	чи болн бидн 'ты и мы'	
$NP \rightarrow Noun (c=x) + Noun (c=x) + Num (f=Conj)$	көвүн күүкн хойр 'мальчик и девочка' (букв. 'мальчик и девочка вдвоем')	
NP -> Noun (c=x) + Noun (c=x)	дегтр, девтр, бичүр 'книга, тетрадь, ручка'	
$NP \rightarrow Pron(c=x) + Noun(c=x)$	би болн эцкм 'я и мой отец'	
NP -> Pron (c=x) + Pron (c=x) + Num (f=Conj)	чи бидн хойр 'ты и я вдвоем' (бук 'ты и мы вдвоем')	

Что касается сочетаний с послелогом, то следует отметить, что можно использовать данные статистики и сведения о грамматике послелогов, о правилах сочетаемости с существительными, стоящими в определенных падежах.

Послелог	Перевод	Падеж
амар	вдоль, у	Gen
	каждый, всякий; еже=; с предыдущим прич. буд.	
_	вр. каждый (всякий) раз, как; по мере того,	3.7
бүр	как, сколько ни, как ни, а, чем, тем	Nom
дараһар	подряд	Nom
деегүр	над, поверх, через, по	Nom
деер	на, возле, около; пока; вместе с, с	Nom; Gen
деерәһәр	поверх, над	Nom
деерәһүр	по	Nom
дор	под, в	Nom
дотр	внутри; в; во; среди, между; в течение	Nom; Gen
дундаһур	сквозь, по центру	Gen
заагар	СКВОЗЬ	Nom
заагур	между, меж, сквозь; не то, не то	Nom; Gen
	с прич. буд. вр. перед, во время, когда; между	
	тем как; при; с указ. мест. и со словами,	
зуур	обозначающими время, переводится по их знач.	Gen
кевәр	согласно, по	Nom
нааһар	почти	Gen
hap	свыше, более	Nom
өмнәһүр	перед, впереди	Gen
өөгүр	вблизи, рядом, мимо	Gen
өөр	у, около, рядом, возле	Gen
сүүләр	после, за	Gen
туршар	в течение, на протяжении, за	Gen
учрар	из-за, ввиду, вследствие; чтоб, чтобы	Gen
хажуһар	около, мимо	Gen
хоорндаһур	между, в промежутке	Gen
шидр	около, вблизи, у, при; к, под	Nom; Gen

Таким образом, рассмотрены структуры синтагм в калмыцком языке, где вершиной является имя существительное. Исходя из рассмотренных синтагм, можно выработать следующие критерии:

- 1) все узлы зависимостей именных групп представляют собой линейную последовательность без разрывов;
 - 2) вершиной являются именные части речи;
- 3) знаки препинания в таких группах отсутствуют, за исключением сочинительных именных последовательностей;
- 4) зависимые элементы находятся слева от вершины, кроме сочиненных и отрицательных комплексов;
- 5) именная синтагма может состоять из одной единицы и более, следовательно, анализатору необходимо обращаться к контексту +1 и более токена.

Границы именных групп могут проводиться в следующих случаях:

- 1) в начале предложения;
- 2) перед знаками препинания;
- 3) после союза хойр 'и';
- 4) перед глаголом в личной форме, а также перед атрибутивными глагольными формами;
 - 5) после отрицательной частицы уга 'нет';
 - 6) перед наречием, если сразу же за ним следует глагол.

Таковы в целом результаты работы по выработке критериев установления границ именных групп, тем не менее работа далека до завершения, возможно уточнение информации о границах и критериях. Следующим этапом станет разработка глагольных групп, при этом уже возможно использовать имеющийся массив данных по именным группам для выработки алгоритма извлечения синтаксических фрагментов из текста.

Литература

Куканова В. В. Принципы семантической разметки Национального корпуса калмыцкого языка // Вестник Калмыцкого института гуманитарных исследований РАН. 2014. № 2. С. 137–143.

Куканова В. В., Каджиев А. Ю. Алгоритм работы морфологического парсера калмыцкого языка // Письменое наследство и информационные технологии [Текст]: Материалы от V Междунар. науч. конф. (Варна, 15—20 сентября $2014 \, \Gamma$.) / отв. ред. В.А. Баранов, В. Желязкова, А.М. Лаврентьев. София; Ижевск, $2014. \, \mathrm{C}. \, 116–119.$

Маслов Ю. С. Введение в языкознание. М.: Высшая школа, 1987. 272 с.

References

Kukanova V. V. Principles of a Semantic Marking of the National Corpus of the Kalmyk Language. *Bulletin of Kalmyk Institute of Humanitarian Research of the RAS*. 2014. No. 2. Pp. 137–143. (In Russ.)

Kukanova V. V., Kadzhiev A. Yu. Algorithm of the Morphological Parser of the Kalmyk Language. In: Written Heritage and Information Technology. Conf. proc. (Varna. 15–20 September 2014). V. A. Baranov, V. Zhelyazkova, A. M. Lavrentiev. Sofia; Izhevsk, 2014. Pp. 116–119. (In Russ.)

Maslov Yu. S. Introduction to Linguistics. Moscow: Vysshaya shkola, 1987. 272 p. (In Russ.)