

## Enhancement Possibilities for the Georgian National Corpus

KAMARAULI MARIAM, *POSTDOC*  
UNIVERSITY OF HAMBURG  
HAMBURG, GERMANY

ORCID: 0009-0006-0404-4424

DOI: [HTTPS://DOI.ORG/10.62343/CJSS.2023.245](https://doi.org/10.62343/CJSS.2023.245)

### ABSTRACT

The aim of this paper is to make suggestions for improving the Georgian National Corpus based on selected linguistic processes. The Georgian National Corpus is currently the most developed and detailed corpus of the Georgian language. One of the reasons for this is the included annotation of the texts, the variety of text genres, and the size of the corpus. While the morphosyntactic analysis of the texts is great, there is room for improvement in the semantic-pragmatic analysis, especially as far as the semantic-pragmatic analysis of functional elements is concerned. Many factors make this issue very interesting, such as grammaticalisation processes or the fundamental development of language. Implementing this type of analysis is essential, especially when it comes to adequate translations by machine translations. The paper contains an approach for analysing functional elements using the example of the particle *xom*.

*Keywords:* Corpus linguistics, Annotation, Modern Georgian, GNC

## **INTRODUCTION**

The 21st century, along with the rapid development of information technologies, brought significant changes to any scientific field and, of course, also to linguistics. The classical grouping of languages established in linguistics has been replaced by a new paradigm of classification. If the traditional classification paradigm included genetic (classification of languages into families according to their genetic relationship), typological (classification of languages according to their morphological structure) and relational classification (classification of languages according to their relational type into, e.g. nominative-accusative, ergative-absolutive and active-stative alignment), today the paradigm of language classification has changed and the focus of language classification added to the quality of the languages' digital representation. What is meant here is the existence of big data both from a quantitative point of view (textbases and speech data of hundreds of millions of tokens) and from a qualitative point of view (high level of annotation quality, electronic dictionaries, grammar resources such as bases of grammatical morphemes and rules, sentiment analysis, treebank, etc.). Thus, according to the approach of language classification, languages are grouped into High Resource Languages (HRL) and Low Resource Languages (LRL). Of the alleged 7,000 languages in the world, only 20 languages have sufficient resources to perform the tasks of Natural Language Processing (NLP). Despite the fact that a large number of monolingual and bilingual digital resources have been created for the Georgian language (GNC, 2024; Georgian Dialect Corpus, 2024; Rustaveli Goes Digital - Parallelkorpus, 2024), it is still classified as a low-resource language (see RichardLitt, 2024). To change this status of the Georgian language, a number of tasks need to be solved, such as the enhancement and further development of the Georgian National Corpus (GNC) – some of the proposals will be presented below.

In general, during the construction of a corpus, the general principles of corpus construction (corpus structure) should be considered, on the one hand, and on the other hand, the structural and grammatical features of the language of the resource embedded in the corpus, which will be considered when creating the corpus search system - the corpus manager. For the efficient use of the corpus, the methodological aspect is also important, in particular, the relationship between data and theory (theoretical qualification of data), the so-called 3A perspective (Wallis & Nelson, 2001: 311ff), namely annotation, abstraction and analysis:

- “Annotation consists of the application of a scheme to texts. Annotations may include structural markup, part-of-speech tagging, parsing, and nu-

merous other representations.

- Abstraction consists of the translation (mapping) of terms in the scheme in a theoretically motivated model or dataset. Abstraction typically includes linguist-directed search but may include rule-learning for parsers, for example.
- Analysis consists of statistically probing, manipulating and generalising from the dataset. Analysis might include statistical evaluations, optimisation of rule-bases or knowledge discovery methods” (Agapova, 2014, p. 282).<sup>1</sup>

The advantage of an annotated corpus is that users can use it for a wider range of research issues and conduct experiments using the corpus manager.

The higher the degree of annotation in the corpus, that is, the more annotation levels are provided in the corpus, the more useful the given corpus is for interdisciplinary research, on the one hand. On the other hand, annotated corpora are needed to implement natural language processing (NLP) and to train artificial intelligence (AI) for a given language.

## ***METHODS***

Two extensive databases have to be mentioned when discussing the Georgian language, namely Thesaurus Indogermanischer Text- und Sprachmaterialien (TITUS) (University of Frankfurt, n.d.) and Georgian National Corpus (GNC) (Georgian National Communications Commission, n.d.). The former comprises corpora of ancient Indo-European languages (such as Avestan, Vedic Sanskrit, Phrygian, or Umbrian) and also materials in more recent Indo-European as well as neighbouring languages, among them the South Caucasian languages (such as Georgian, Megrelian, Svan and Laz) but TITUS does not contain as many textual resources for Modern Georgian as GNC. The National Corpus of the Georgian Language (GNC) is the largest corpus created for the Georgian language (more than 202 million tokens), which is the reason. GNC belongs to the type of diachronic corpora, which com-

<sup>1</sup> Wallis, S. (n.d.). Annotation takes a set of texts and adds linguistic information to it, enriching it and identifying instances of linguistically meaningful entities and relations. At this point, the resulting enriched dataset (‘corpus’) is usually distributed to the research community. Abstraction is the researcher’s exploratory process of establishing a mapping between concepts they wish to research, and representations found in the corpus (text + annotation). It also maps the structured corpus to a regular dataset that can be analysed by conventional statistical methods. The key linking element in abstraction is a corpus query. Analysis is the process of applying statistical and other methods to data that has been abstracted in this way. Retrieved from <https://www.ucl.ac.uk/english-usage/staff/sean/>

bines Old, Middle, and Modern Georgian language resources. The corpus includes both resources of the written Georgian language from ancient monuments (inscriptions, handwritten sources) to the present day, and samples of oral speech - the Georgian dialect corpus is integrated into the corpus. When it comes to text genres, GNC is a balanced corpus containing religious, historical, juridical and political texts. The latter two genres are also represented as separate sub-corpora. Nevertheless, the corpus requires further development both in terms of genre and quantity.

## Figure 1

### *The sub-corpora of the GNC*

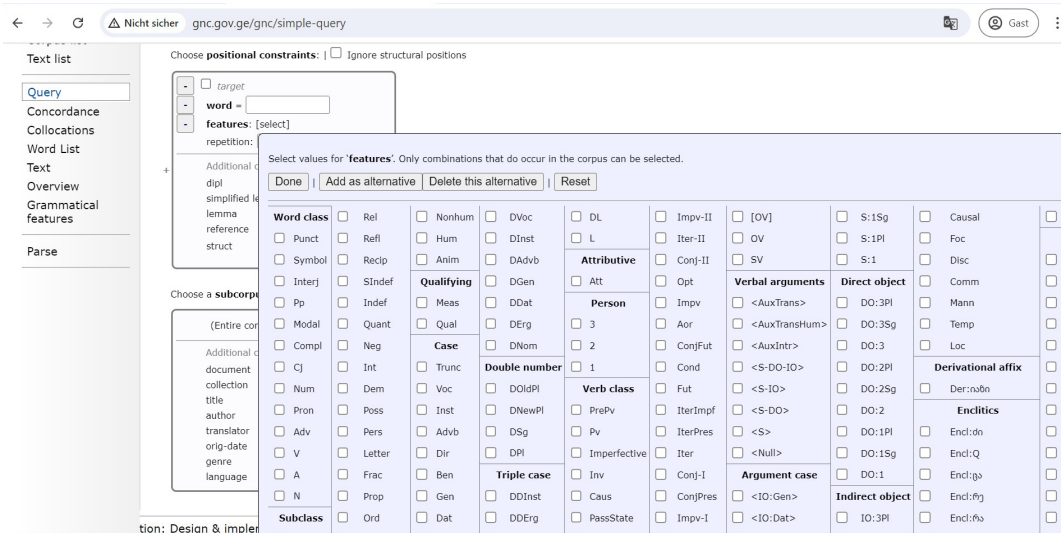
The screenshot shows the GNC website interface. At the top, there is a navigation bar with the GNC logo and the text 'ქართული ენის ეროვნული კორპუსი' and 'The Georgian National Corpus'. Below this, there is a sidebar with a menu containing 'GNC Home', 'About the project', 'Using the GNC', 'Documentation', 'Publications', 'Corpus list' (highlighted), 'Text list', 'Query', 'Concordance', 'Collocations', 'Word List', 'Text', 'Overview', 'Grammatical features', and 'Parse'. The main content area is titled 'Corpus list' and contains the following table:

Corpus	Size (words & punctuation)	Updated	Description
<b>GNC Old Georgian</b>	7 101 021	2022-12-31	Georgian National Corpus, Old Georgian
<b>GNC Middle Georgian</b>	1 432 262	2019-06-25	Georgian National Corpus, Middle Georgian
<b>GNC Modern Georgian</b>	1 993 022	2023-01-01	Georgian National Corpus, Modern Georgian
<b>GRC</b>	202 728 329	2016-12-05	Georgian Reference Corpus
<b>GDC</b>	1 694 362	2015-09-14	Georgian dialect corpus
<b>GNC Political texts</b>	1 436 075	2019-08-06	Georgian National Corpus, Political texts
<b>GNC Law texts</b>	1 495 985	2019-04-15	Georgian National Corpus, Old and Middle Georgian, Law texts
<b>GNC Megrelian</b>	89 404	2015-09-14	Georgian National Corpus, Megrelian
<b>GNC Svan</b>	473 180	2015-09-14	Georgian National Corpus, Svan

In addition to the Georgian language, the GNC includes resources for other South-Caucasian languages - Megrelian and Svan. Both the textual material published in these languages and the modern oral resources (which only represent a fraction of what TITUS has to offer) were obtained and processed within the framework of the international scientific projects implemented at the University of Frankfurt (TITUS, ECLinG, SSGG), are presented here. A large Georgian reference corpus (GRC) is included, which contains less thoroughly processed texts from various fictional and non-fictional domains.

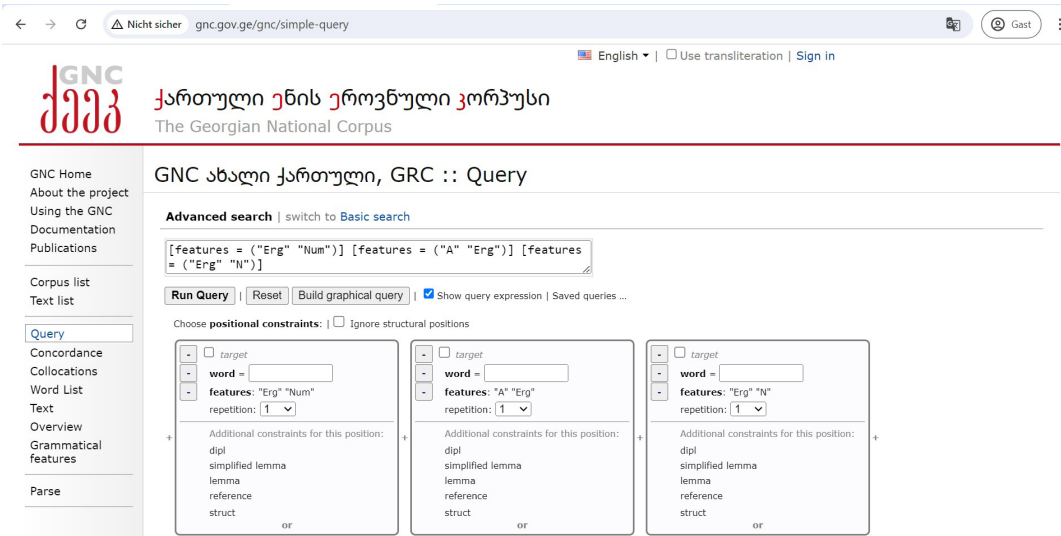
GNC is an annotated corpus - the corpus manager allows for both simple and complex searches in the corpus. In the case of a complex search, it is possible to find a word form according to one or several grammatical features combined.

Figure 2  
Example of a complex search in the GNC



The corpus search engine also allows you to search the corpus for phrasal constructions:

Figure 3  
Searching interface for phrasal constructions in the GNC (a phrase containing a numeral, an adjective and a noun in the ergative)



The results of the search are then displayed in the corresponding concordance:

## Figure 4 Results of the search

corpus	cpo	match \.
gnc-kat	1329646	უც და ახლა წელსვეთი მიშველი დავჯეკი, წვიძის ათსმა, ათი ათსმა, ასი ათსმა, ასი ათსმა, მხამაღ! მხამაღ! </p><p> </b> თეიმურაზს ერთ წუთში თვალწინ
gnc-kat	1406893	სახარე </b> ვაჰარი გიორგი ერისთავმა, ცაგარელმა, ანტიონოვმა </b> და
gnc-kat	1318913	ი კომლექსი არაფერს ხიზნავს. პოეტს სულში უნდა ჩახვდეთ, – თქვა ერთ-
gnc-kat	611206	და ვწნათ ამ უგზო ქვეყანაში? </p><p> </b> – ვაი ჩემს გამჩეხს! – უკასუსხა
gnc-kat	418303	ყ კანმა რომ თუკას, რატომ უნდა გვერდითი?! – დავებრა მოულოდნელად. –
gnc-kat	1086533	ერთმა თქვენიდანაა გლეხმა
gnc-kat	670550	ერთმა კეთილმა დისახლისმა
gnc-kat	906857	ერთმა მოხა კალმა
gnc-kat	609849	ერთმა მისხველიმა ყმაწვილ-კავამ
gnc-kat	406999	ერთმა მოხუცებულმა სოფელელმა
gnc-kat	2044127	ერთმა მსუქანმა კალმა
gnc-kat	976205	ერთმა მშვენიერმა დედმა
gnc-kat	976197	ერთმა მშვენიერმა დედმა
gnc-kat	2100699	ერთმა პარხელმა ქართველმა
gnc-kat	1211444	ერთმა სახარეო მანქანამ
gnc-kat	568491	ერთმა უბრალო მშობვევამ
gnc-kat	433965	ერთმა უგზურმა თავუნამ
gnc-kat	1026825	ერთმა უსამართლო მშობელიამ

The high degree of annotation in the corpus allows for morphosyntactic and syntactic analysis:

## Figure 5 Parsing of a sentence

Parse

Here you can parse sentences with the morphological analysers that are being used to analyze the texts of the GNC. Write a sentence or shorter text and click the 'Parse' button. You can click on the morphosyntactic features of the analyzed text to have them explained.

Language variety: Modern Georgian

მე ძლიერ მიყვარს ისევერ თოვლის ქალწულბივითი ხიდიდან ფენა.

Parse Reload Show all readings Show used rules Show dependencies UD features

მე	3	მე	Pron Pers 1 Dat Sg
ძლიერ	4	ძლიერ	Adv Deg
მიყვარს	3	სიყვარულ-ი/ყვარ	V MedPass Inv Pres OV <S:DO> <S:Dat> <DO:Nom> S:1Sg DO:3
ისევერ	6	ისევე[ერ]ი	A Gen Att <OldPl>
თოვლის	1	თოვლ-ი	N Gen Sg
ქალწულბივითი	1	ქალწულ-ი	N Hum Nom Pl NewPl PP PP:3non
ხიდიდან	3	ხიდი-ი	N Inst Sg PP PP:დან
ფენა	6	ფენა/ა/ფენ	N VN Nom Sg
.	1	.	Punct Period

GNC was created within the framework of international scientific cooperation in the years 2012-2019. Both European (Frankfurt University, University of Bergen) and Georgian scientific and educational institutions (Georgian National Communications Commission, n.d.) participated in its creation.

The quality of big data annotation is crucial for AI tasks. The quality of data annotation refers to the accuracy and consistency of data labelling for machine learning models. It is crucial to ensure that the algorithms learn effectively from the annotated data provided. High-quality data annotation leads to more accurate predictions and better model performance. It also implies a multi-level system of analysis, which includes morphological, morphosyntactic, syntactic, pragmatic, and semantic levels. In the case of speech data, in addition to text, audio and video resources, suprasegmental analysis is also provided. Suprasegmental features help to convey meaning, structure and emotional undertones in oral communication. They affect the way syllables, words and sentences are pronounced and influence the meaning and perception of spoken language at a higher level.

The GNC is characterised by a relatively high level of token annotation, which includes both the lemma and grammatical features of the token, as well as other relevant information (source, author, title, date of the text, suprasegmental annotations, etc.). Below, an example from the nominal morphology is provided:

**Figure 6**  
*Search result of the noun მღაზიებში “in the stores”*

The screenshot shows the GNC search interface with the following details:

- Search Interface:** Includes a search bar with the query 'მღაზიებში', a 'Run Query' button, and options for 'Refine', 'window: document', 'Stop', and 'Saved queries...'. It also shows 'Done. Running time: 0.10 sec. (0.03 CPU sec.)' and 'Page size: 500px'.
- Navigation:** 'Hit 1 - 30 of 1301' with 'Previous', 'Next', and 'Go to:' options. There are also 'Download (Excel mode)' and 'Copy query URL' buttons.
- Search Results Table:**

corpus	cpos	match
gnc-kat	696932	მღაზიებში უმეტეოდ რომ დაიქვხნ თურქე საფურას ვიწმე ქალს.
gnc-kat	778242	მღაზიებში შექყიჭა, გადაამკურველისთვის რადაც დაეკალიბინა,
gnc-kat	1379888	იების ბედის გასაგებად და ორიოდე გრომის მისაღებად. </p> <p> </b> </b> იმ
gnc-kat	1379934	ა საყუთარ </b> ცხოვრებისა და ბედის წიგნს. </p> <p> </b> რა არ იყო იმ
gnc-kat	1380097	ი და სუველიანი დასარული. </p> <p> </b> ხანგამომწვევითი თეიმურაზი იმ
gnc-kat	1380205	ე უსაქმური თეიმურაზი ზოგჯერ მთელ დღეს ტრიალებდა </b> საკომისიო
gnc-kat	1787449	ა ასორტიმენტის, მაგრამ უარგისი საკომისიო იქმებოდა, რაც ოფურჩების
gnc-kat	1787474	ყოფად მიიხიშეს და გადაწყვიტეს, უარი ეთქვათ, რათა ეს უარი სხვა დროს
gnc-kat	1788332	ა ზომის, თეთრი, ოვგოსლაკური, ქალის ჩემა ცენტრსა და მის მიმდებარე
gnc-kat	1821612	იხდა თურმე, იქ გაიხარდა. სკოლა რომ დაამთავრა, მისი კლასის გოგონები
gnc-kat	1971398	იხდესო, ეს ფერი ხაზები ძალიან მიხდებაო, ოღონდ სადმე ვიშოვიდეთ. აქ
gnc	63787	ყოფითი რაიმე გაგონ. ისინი მხოლოდ მზამზარეულ საგნებს ციფულობენ
gnc	537720	ა მოვების ბედის გასაგებად და ორიოდე გრომის მისაღებად. </p> <p> </b> იმ
gnc	537758	იხოვლობდა საყუთარ ცხოვრებისა და ბედის წიგნს. </p> <p> </b> რა არ იყო იმ
gnc	537906	ანყისი და სუველიანი დასარული. </p> <p> </b> ხანგამომწვევითი თეიმურაზი იმ
gnc	537991	ა უკვე უსაქმური თეიმურაზი ზოგჯერ მთელ დღეს ტრიალებდა საკომისიო
gnc	1654143	ჭერსაებო? – რის შემდეგაც დარჩენვილინი მივდივართ კარებისაყნ. ასეთ
gnc	1670044	ფიგიაბტი არ დაკითხუვბა. ვითომ არ ვიცი, რატომ დადიხარ ტანსაცმლის
- Match Details:**
  - word: მღაზიებში
  - dipi: მღაზიებში
  - simplified lemma: მღაზია
  - lemma: მღაზი[ა]
  - features: N Dat Pl NewPI PP PP;8n
  - document: NG/chiladze-o/chiladze-o+godori
  - title: გოდორი
  - reference: გოდორი პირველი მაწილი I 35
  - author: ვილაძე, ოთარ
  - genre: /fiction/
  - language: kat

The same applies to search results from the verbal morphology and uninflectable words:

**Figure 7**  
*Search result of the verb ტრიალებდა “[(s)he] was spinning”*

The screenshot shows the GNC concordance search interface. The search term is 'ტრიალებდა'. The results table has columns for 'corpus', 'cpos', and 'match'. The 'match' column shows the word 'ტრიალებდა' and its grammatical features: 'word: ტრიალებდა', 'dipl: ტრიალებდა', 'simplified lemma: ტრიალი', 'lemma: ტრიალ/იტრიალ', 'features: V MedAct Impf <S> <S:Nom> S:3Sg', 'document: NG/amiredzhibi-ch/amiredzhibi-ch+data-tutashxia', 'title: დათა თეთაშხია', 'reference: 8', 'author: ამირჯიბი, ჭახუა', 'orig-date: 1974', 'genre: /fiction/', 'language: kat'.

**Figure 8**  
*Search result of the affirmative particle xom (ხომ)*

The screenshot shows the GNC concordance search interface. The search term is 'ხომ'. The results table has columns for 'corpus', 'cpos', and 'match'. The 'match' column shows the word 'ხომ' and its grammatical features: 'word: ხომ', 'dipl: ხომ', 'simplified lemma: ხომ', 'lemma: ხომ', 'features: Adv Disc', 'document: NG/mishveladze-r/mishveladze-r+tom1-04', 'title: რჩეული თხზულებანი IV - მოველები', 'reference: ქველი 169', 'author: მიშველაძე, რევაზ', 'genre: /fiction/', 'language: kat'.

In the case of uninflectable words, as shown in Fig. 8, the syntactic-pragmatic function is indicated: *xom* - Adv Disc (discourse adverb). However, a certain part of the tokens in GNC is not annotated, which is due to the fact that the issues of the functional grammar of the Georgian language are still theoretically unresearched and have only been studied in fragments. Accordingly, the grammatical character-



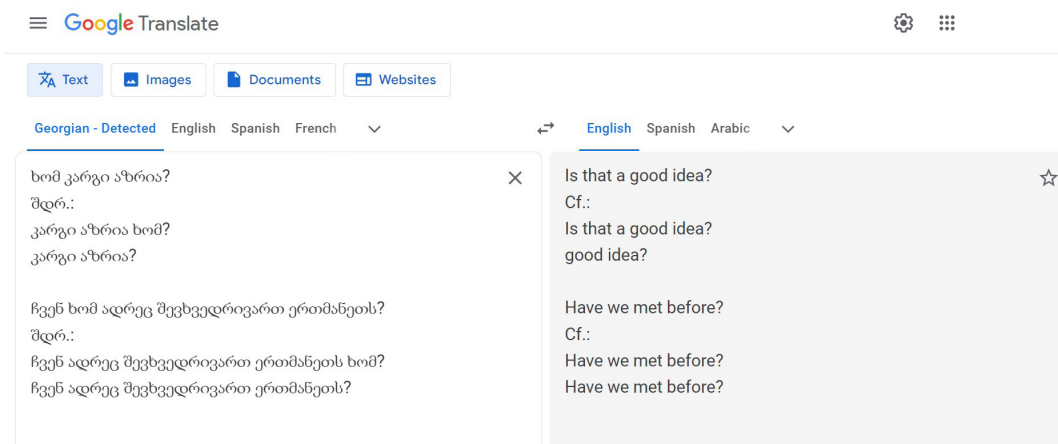
sation of some words in the corpus is either inaccurate or the grammatical features are not defined at all - in such a case, only “unknown” is indicated. Below, we discuss several works related to GNC annotation system improvement and corpus development proposals and present my proposal regarding the annotation of invariant words in the corpus.

## RESULTS

The functional-semantic analysis of uninflectable elements is utterly significant not only for refining the theoretical model and for creating a functional grammar of the Georgian language, but it also has practical significance for the development of language technologies, especially for the improvement of automatic translation. None of the currently available translation programs can adequately convey the semantic difference in Georgian sentences from a functional-semantic point of view:

### Figure 10

*How Google translates sentences with and without the particle *xom**



As shown in Fig. 10, Google Translate does not differentiate between the meaning of sentences with or without the particle *xom*, which makes the significance of such an analysis all the more necessary.

## DISCUSSION

In order to achieve a high-quality annotation, specific phenomena of any given language must be considered - structural features, grammatical processes in the language, functional-semantic and pragmatic meaning of linguistic elements, and

other specific features. This applies not only to simple elements of the corpus, such as word forms, but also to complex structural units, such as phrasal structures. In my opinion, further development of GNC requires the refinement of the specific phenomena of the Georgian language. Below, I will present suggestions for improving the annotation of the Georgian National Corpus, using examples of simple elements and complex constructions.

One of the linguistic phenomena of the Georgian language is approximative verbs; these elements represent a symbiosis of the nominal and verbal domain, as the marker used for approximateness originates from the nominal domain and is suffixed to a fully inflected verb. The suffix *-vit* ('like, as'), which is typically suffixed to a noun in the nominative or the dative case (the former applies to nouns with consonantal stems, the latter to nouns with vocalic stems), can also be found suffixed to nouns in the genitive case, which is the rarest among the cases in combination with the suffix (Kamarauli, 2023, p.52).

## Figure 9

Example of an approximative verb, classified as “unknown”

The screenshot shows the GNC Concordance interface. The search results table is as follows:

corpus	cpos	match
grc	16856225	მომიბოდიმასავით ქალბატონმა მაროკამ. </p> <p> იმდენად უა'
grc	100499152	დური თუა, გერმანიაშიც ვიყიდი და ტყუილად ჩაბთა უხდა დავიძიძიოო - მკვების ასეთ მამდილო მიზ
grc	180292888	წერითი ოჯახი როგორ ვარჩიზი, გლუხის შვილი მინსაც თუ არ ვერჩივო, - ქმედებაზე ჩამოვარდა. ...
grc	182374743	ინით ვერასოდეს ვუკავშირდებოდი. როდესაც კუჩაში შემთხვევით შეხვდა, , რაკი მაკოორიტარი დეპუ

The detailed view for the first hit shows the following metadata:

- word: მომიბოდიმასავით
- norm: მომიბოდიმასავით
- lemma: -
- features: Unknown
- source: lib.ge
- author: ბიზო ჯიბლაძე
- title: სიყვარული ამერყამი

The GNC has not yet provided a classification for such constructions, so these are labelled as “unknown”. What I propose is the following: when verbs are analysed as usual according to the grammatical markers such as person, number, tense, etc., another feature must be added, namely verbal approximateness (AppV). The morpheme expressing approximateness (APP: სავით) should be added at the end of the grammatical features:

Cf.:

EXAMPLE	GRAMMATICAL FEATURES
მეცადინეობდა	V MedAct Impf <S> <S:Nom> S:3Sg
vs.	
მეცადინეობდასავით	AppV MedAct Impf <S> <S:Nom> S:3Sg APP: სავით

One of the important challenges in the analysis of the Georgian language is the issue of annotation of uninflectable elements - particles, conjunctions, adverbs, conjunctions. The correct annotation of functional elements is indispensable for solving both semantic analysis and treebank tasks.

Below, I present my annotation approach of functional elements on the example of the functional-semantic analysis of the particle *xom*.

The particle *xom* is analysed as an interrogative particle in scientific literature, in particular as:

- An interrogative particle, which 1. is used in interrogative clauses and denotes confirmation, and 2. is used together with a negative word (არ, ვერ, არავინ...) and indicates doubt (Explanatory Dictionary n.d);
- An interrogative particle-morphemoid, which a) expresses confirmation in interrogative clauses, b) expresses doubt with negative morphemoids (no, can, nobody), c) is used in negative constructions to express the function of a request (Jorbenadze, K'obakhidze & Beridze, 1988: 474-475);
- It is used when asking a question and wanting to have the answer confirmed (Georgian Dictionary n.d);
- It is annotated as a discourse adverb in the National Corpus of the Georgian language (Georgian National Corpus n.d.).

In the reference sub-corpus of GRC, *xom* is statistically one of the most frequently used particles. Table 1 (see next page)

The functional-semantic analysis of the particle *xom*, which is presented below, relies on the resources provided by the GNC. Both classic research methods and corpus linguistic research methods are used to analyse the examples. Additionally, substitution, elimination, permutation and paraphrasing tests were also used in the research. The corpus linguistic analysis showed that the particle can convey more functional semantics than in the definitions presented above. In addition, the conducted analysis showed that the following parameters are crucial for determining the functional semantics of the particle *xom*, which will be introduced below:



As the examples above show, it is possible to transpose the particle *xom* in (1a-b) and even omit (1c) from the sentence. In the case of transposition, the sentence maintains the semantics of confirmation (affirmativeness). Therefore, the probable answer is ‘yes’. In the case of omission, affirmativeness is lost, and the sentence becomes a ‘yes/no’ question - the answer can be either positive or negative.

Both sentences (1a) and (1b) require a positive answer. The difference between them is the speaker’s attitude: in (1a), the speaker offers his opinion to the listener, which is affirmative and conveys the speaker’s position; as a result of the transposition of the particle in (1b), the speaker expects the listener to confirm the opinion expressed by him.

The following example confirms that the particle *xom* placed in the final position expresses the expectation of confirmation from the listener:

(2a)	<i>ramden-ze</i>	<i>gagvarige</i>	<i>me</i>	<i>da</i>	<i>besarion-i?</i>
	how much.DAT.SG-on	settle.S2SG.O1PL.AOR	I.NOM.SG	and	Besarion-NOM.SG
	<i>otxas-i</i>	<i>manet-i</i>	<i>unda</i>	<i>moeca</i>	<b><i>xom?</i></b>
	fourhundred-NOM.SG	Mane-ti-NOM.SG	MPTCL	give.S3SG.PLUPERF	AFF

‘How much money did me and Besarion agree on thanks your help? He should have given me 400 Manetis, right?’ (Davit kldiašvili, *Soloman Morbelaze*)

When the particle *xom* is placed in the initial position, the speaker expects the listener to confirm the amount of money:

(2b)	<i>ramden-ze</i>	<i>gagvarige</i>	<i>me</i>	<i>da</i>	<i>besarion-i?</i>
	how much.DAT.SG-on	settle.S2SG.O1PL.AOR	I.NOM.SG	and	Besarion-NOM.SG
	<b><i>xom</i></b>	<i>otxas-i</i>	<i>manet-i</i>	<i>unda</i>	<i>moeca?</i>
	AFF	fourhundred-NOM.SG	Maneti-NOM.SG	MPTCL	give.S3SG.PLUPERF

‘How much money did me and Besarion agree on thanks your help? He should have given me 400 Manetis, right?’

Example (2a) is an interrogative clause, and the answer requires specifying the amount. In the following example, (2b), the speaker states the amount himself and waits for the addressee to confirm it. Both sentences are affirmative sentences, but



In this context, the particle *xom* is a pragmatic element, namely a presupposition marker. If we omit the adversative conjunction *magram* ‘but’ in the last sentence, we get the following expression: *tkven xom geograpi xart?* ‘You are a geographer, **right?**’. Here, the presupposition is clearly readable, and it is marked in the sentence with the particle *xom*. By eliminating it, the presupposition in the sentence is lost - the sentence turns into a simple ‘yes/no’ question: *tkven geograpi xart?* ‘Are you a geographer?’. The adversative conjunction *magram* ‘but’ makes the speaker’s position even stronger: the geographer’s answers in the discourse (lack of geographical knowledge) surprise the speaker since he expects the geographer to have this knowledge. The opinion of the speaker in the last sentence is critical, which is marked by the adversative conjunction *magram* in the initial position, and to convey his position, the speaker uses an affirmative sentence with the particle *xom*.

- Final position and scope effects

The possibility to transpose elements also brings some changes in scope and, therefore, semantics. The following examples have been constructed to demonstrate the functionality and the resulting scope effects of the particle *xom* when transposed:

(4a) *xom*            *luḡa-m*            *dalia*            *sam-i*            *lud-i?*  
 AFF            Luka-ERG.SG            drink.S3SG.AOR            three-NOM.SG            beer-NOM.SG

‘Luka drank three beers, right?’

(4b) *luḡa-m*            *xom*            *dalia*            *sam-i*            *lud-i?*  
 Luka-ERG.SG            AFF            drink.S3SG.AOR            three-NOM.SG            beer-NOM.SG

‘Luka drank three beers, right?’

(4c) *luḡa-m*            *dalia*            *xom*            *sam-i*            *lud-i?*  
 Luka-ERG.SG            drink.S3SG.AOR            AFF            three-NOM.SG            beer-NOM.SG

‘Luka drank three beers, right?’

\*(4d) *luḡa-m*            *dalia*            *sam-i*            *xom*            *lud-i?*  
 Luka-ERG.SG            drink.S3SG.AOR            three-NOM.SG            AFF            beer-NOM.SG

‘Luka drank three beers, right?’

(4e) *luḡa-m*            *dalia*            *sam-i*            *lud-i*            *xom?*  
 Luka-ERG.SG            drink.S3SG.AOR            three-NOM.SG            beer-NOM.SG            AFF

‘Luka drank three beers, right?’





The opinion that such sentences serve as argumentations is methodologically difficult to justify in the case of simple sentences, but in case of more complex syntactic constructions, we can the method of paraphrasing:

- (7a) *çarmodgena-c ara akvs mosalodnel saprtxe-ze,*  
 idea.NOM.SG-FOC NEG have.S3SG.PRES expecting.DAT.SG danger.DAT.SG-on
- gavipikre me. mas xom arasodes gamoucdia*  
 think.S3SG. I.NOM.SG he.NOM. AFF never experience.S3SG.PERF  
 AOR SG
- šimšil-i da çqurvil-i*  
 hun- and thirst-NOM.SG  
 ger-NOM.SG

‘He has no idea about the impending danger, I thought. - He has never experienced hunger and thirst.’ (Antoine de Saint-Exupéry, *The Little Prince*)

→ Paraphrasing the second sentence

- (7b) *vinaidan mas araso- gamoucdia šimšil-i*  
 as he.NOM.SG never experience.S3SG.PERF hunger-NOM.  
 SG
- da çqurvil-i*  
 and thirst-NOM.SG

‘As he has never experienced hunger and thirst.’

- (8a) *me unda vizruno mas-ze. igi xom*  
 I.NOM.SG M P T- care.S1SG.OPT ( s ) h e . ( s ) h e . N O M . S G AFF  
 CL DAT.SG
- iset-i sust-i da iset-i gulubrçqvil-a*  
 such-NOM.SG weak-NOM.SG and such-NOM.SG naïve.NOM.  
 SG-COP

‘I have to care about her/him. He is so weak and so naïve.’ (Antoine de Saint-Exupéry, *The Little Prince*)

→ Paraphrasing the second sentence

- (8b) *vinaidan igi iset-i sust-i da*  
 as (s)he.NOM.SG such-NOM.SG weak-NOM.SG and
- iset-i gulubrçqvil-a*  
 such-NOM. naïve.NOM.  
 SG SG-COP

‘As he is so weak and so naïve.’

As shown in the examples (7a-b) and (8a-b), we can consider that the particle *xom* is used as an argumentation marker when it is realised in the midfield of declarative sentences.

In interrogative sentences, the particle can be realised in combination with the modal word *šeizleba* ‘can’ (1635 such cases are confirmed in the GNC) and conveys possibility, permission or assumption in all three positions:

- (9) [...] ***xom*** *šeizleba* *tan* *rağac* *gkitxot?*  
 [...] AFF can at the something.NOM.SG ask.S1SG.O2PL.OPT  
 same time

‘[...] I can ask you something at the same time, right?’ (Journal *Liṭeraṭuruli p̄aliṭra*, 2008)

- (10) *magram kac-i-c* ***xom*** *šeizleba* *iqos* *meçq̄vile!*  
 but man-NOM.SG-FOC AFF can be.S3SG. partner.NOM.SG  
 OPT

‘But a man can also be a partner, can’t he!’ (Tariel Čanṭuria, *Orni kuṭeši*)

- (11) *šen-tan ertad rom ṭrailer-it vimgzavro, xom šeizleba?*  
 you.DAT.SG- together that trailer-INST. travel.S1SG. AFF can  
 WITH SG OPT

‘Is it possible for me to travel with you in a trailer?’ (Aḳaḳi Gegenava, *Mogzauris d̄giurebi*)

The combination *xom šeizleba* can also be used in declarative clauses:

- (12a) *magram zogžer vpikrob: xom šeizleba rom*  
 but sometimes think.S1SG.PRES AFF can that  
  
*adamian-s sakme daavicq̄des.*  
 human-DAT. business. forget.S3SG.O3SG.OPT  
 SG NOM.SG

‘But sometimes I think: a human can forget about the business, can’t he.’ (Antoine de Saint-Exupéry, *The Little Prince*)

- (12b) *magram zogžer vpikrob: šeizleba xom rom*  
 but sometimes think.S1SG.PRES can AFF that  
  
*adamian-s sakme daavicq̄des?*  
 human-DAT. business. forget.S3SG.O3SG.OPT  
 SG NOM.SG

‘But sometimes I think: a human can forget about the business, right?’

(12c)	<i>magram</i>	<i>zogžer</i>	<i>vpikrob:</i>	<i>šeižleba</i>	<i>rom</i>	<i>adamian-s</i>
	but	sometimes	think.S1SG.PRES	can	that	human-DAT.SG
	<i>sakme</i>	<i>daavicqdes</i>				
	business.	forget.S3SG.O3SG.OPT				
	NOM.SG					

‘But sometimes I think: can a human forget about the business?’

In the case of elimination of the particle *xom* as shown in (12c), the dependent clause requires a transformation into an interrogative clause, which can function as a rhetorical question. The paraphrase of this sentence would be: “A man cannot forget his work.” In the case of the transposition of the particle *xom* in the second position in (12b), the affirmative sentence with the semantics of possibility is preserved, but the perspective changes: the speaker expects to receive confirmation from the listener.

The combination *xom šeižleba* can also be in the second position as in the next example, and here too the particle *xom* conveys the expectation of the speaker to receive confirmation:

(13)	<i>ese-c</i>	<i>xom</i>	<i>šeižleba</i>	<i>iqos</i>	<i>liṭeraṭura?</i>
	this.NOM.SG-FOC	AFF	can	be.S3SG.OPT	literature.NOM.SG
	<i>aman-a-c</i>	<i>xom</i>	<i>šeižleba</i>	<i>bestseler-is</i>	<i>saxel-i</i>
	this.ERG.SG-EMPH.V-FOC	AFF	can	bestseller-GEN.SG	name-NOM.SG
	<i>moixvečos?</i>				
	gain.S3SG.O3SG.OPT				

‘This can also be literature, right? This can also gain the title of a bestseller?’ (Nene Kṽiniḳaḻe, *Iaguarebis tekno*)

Declarative clauses with the particle *xom* are characterised by more intensity, the persuasive power of the opinion expressed by the speaker is greater, which is strengthened by the repetition method used in this case. Accordingly, these types of sentences are often found in the speeches of politicians.

In interrogative sentences, the particle *xom* is often found in combination with the negation particle *ar*, although the negation particle itself is not desemanticised (also called semantic bleaching), but the meaning of the sentence does not convey negation on a pragmatic level. In such sentences, both particles *xom* and *ar* should be considered as one functional element ‘*xom+ar*’. In case of transposition and elim-

ination, they are moved or eliminated together. The combination of *xom* and *ar* is used in the initial position during a polite question:

- (14a) *xom ar gciva?* vs. (14b) *gciva?*  
 AFF NEG being cold.S2SG.PRES being cold.S2SG.PRES  
 ‘You are not feeling cold, are you?’ ‘Are you cold?’
- (15a) *xom ar dagaviçqdeba?* vs. (15b) *dagaviçqdeba?*  
 AFF NEG forget.S2SG.O3SG.FUT forget.S2SG.O3SG.FUT  
 ‘You won’t forget, will you?’ ‘Will you forget?’
- (16a) *xom ar geçqineba?* vs. (16b) *geçqineba?*  
 AFF NEG being offended.S2SG.O3SG.FUT being offended.S2SG.O3SG.FUT  
 ‘You won’t feel offended, will you?’ ‘Will you feel offended?’

The combination *xom+ar* is mostly found in the second position, and depending on which verb it is combined with, it conveys different semantics:

- Questions with propositional semantics:

- (17a) *rame xom ar ginda?*  
 something.NOM.SG AFF NEG want.S2SG.O3SG.PRES  
 ‘Do you want anything?’ (Akaçi Gegenava, *Mogzauris dği-urebi*) → offering to bring/buy
- vs.
- (17b) *rame ginda?*  
 something.NOM.SG want.S2SG.O3SG.PRES  
 ‘Do you want something?’ → yes/no question

- Question with the semantics of doubt:

- (18a) *brma xom ar aris?*  
 blind.NOM.SG AFF NEG be.S3SG.PRES  
 ‘(S)he isn’t blind, is (s)he?’ (Niço Lomouri, *Paçia megobrebi*) → expressing doubt
- vs.
- (18b) *brma aris?*  
 blind.NOM.SG be.S3SG.PRES  
 ‘Is (s)he blind?’ → yes/no question



- In initial or final position, the particle *xom* refers to the whole sentence but triggers different readings:
  - a. In initial position, the affirmation requires confirmation from the perspective of the listener;
  - b. In final position, the affirmation is given from the perspective of the speaker;
- The particle *xom* refers to entire phrases and not to single elements of phrases;
- When combined with the negation particle *ar*, the combination *xom+ar* has to be considered one functional element;
- Depending on the position *xom+ar*, the sentence can have different semantics:
  - a. In initial position, the sentence can convey politeness;
  - b. In midfield position, the following semantics can be conveyed:
    - i. Propositional semantics,
    - ii. Semantics of doubt,
    - iii. Clarification/inquiring semantics,
    - iv. Rhetorical question.

The analysis of the particle *xom* showed that in order to accurately understand and translate Georgian, not only a morphosyntactic but additionally a semantic-pragmatic analysis should be implemented. Of course, there are still many relevant aspects left to research; this paper served to present a first approach and to open the topic for future research.

**ABBREVIATIONS**

ADV	adverbial case	MPTCL	modal particle
AFF	affirmative	NEG	negation
AOR	aorist tense/aspect	NOM	nominative case
COP	copula	O	object
DAT	dative case	OPT	optative
EMPH.V	emphatic vowel	PERF	perfect tense/aspect
ERG	ergative case	PLUPERF	plusquamperfect
EXT.V	extensional vowel	PRES	present tense
FOC	focus	PL	plural
FUT	future tense	S	subject
GEN	genitive case	SG	singular
INST	instrumental case	1/2/3	1 <sup>st</sup> /2 <sup>nd</sup> /3 <sup>rd</sup> person

**REFERENCES**

Agapova, S. (2014, June). On text linguistics. In *Collected Articles of the 3rd International Linguistics Conference* (Taganrog, Russia) (p. 265). Cambridge Scholars Publishing.

Association of language modelling. (n.d.). Enis modelirebis asociacia. Retrieved on June 15, 2024, from <https://www.ena.ge/explanatory-online>

Big Georgian dialect base and the electronic atlas of Georgian dialects GEDA. (n.d.). *Didi kartuli dialect'uri baza da kartuli dialekt'ebis elekt'ronuli at'las GEDA*. Retrieved on June 15, 2024, from <http://www.corpora.co>

Georgian Lexicon. (n.d.). *Kartuli leksik'oni*. Retrieved on June 16, 2024, from <https://www.ganmarteba.ge/word/ბმმ>

GNC. (2024, June 15). *Georgian Dialect Corpus*. K'art'uli enis erovnuli k'orp'usi. Retrieved on June 15, 2024 from <http://gnc.gov.ge/>

Jorbenadze, B., K'obakhidze, M., & Beridze, M. (1988). *Kartuli enis morpemebisa da modaluri element'ebis leksik'oni (Masalebi kartuli enis sist'emati'uri k'ursistvis)* [The lexicon of morphemes and modality elements of the Georgian language (Collection for the systematic course of the Georgian language)]. Tbilisi: Publishing House "Metsniereba".

Kamarauli, M. (2023). Approximative verbs: A symbiosis of the nominal and the verbal domain (on the example of the Georgian language). *Millennium*, 1, 50-70.

Low resource languages. (n.d.). Retrieved on June 15, 2024, from <https://github.com/RichardLitt/low-resource-languages>

TITUS - Thesaurus Indogermanischer Text- und Sprachmaterialien. (n.d.). Retrieved on June 15, 2024, from <https://titus.fkidg1.uni-frankfurt.de/>

Wallis, S., & Nelson, G. (2001). Knowledge discovery in grammatically analysed corpora. *Data Mining and Knowledge Discovery*, 5, 305-335.

Rustaveli Goes Digital - Parallelkorpus. (n.d.). *Parallel Corpus of the Translations of the epic "The Knight in the Panther's Skin"*. Retrieved on June 15, 2024, from <https://rustaveli-goes-digital.de>