#### Paul Meurer

AKSIS, University of Bergen, Norway

Seventh International Tbilisi Symposium on Language, Logic and Computation Tbilisi, October 1 – 5, 2007

### **Outline**

- 1 Introduction: The Georgian grammar project
- Lexical-Functional Grammar
- Morphology and Morphosyntax
- Some aspects of Georgian syntax
- Tools for LFG grammar development

#### **Outline**

- 1 Introduction: The Georgian grammar project
- Lexical-Functional Grammar
- Morphology and Morphosyntax
- Some aspects of Georgian syntax
- 5 Tools for LFG grammar development

#### The Georgian grammar project

 is aimed at developing a full-scale computational grammar for Georgian

- is aimed at developing a full-scale computational grammar for Georgian
- uses LFG (Lexical Functional Grammar) as a syntactic formalism (output of a parse consists of c- and f-structures, no semantic representation yet)

- is aimed at developing a full-scale computational grammar for Georgian
- uses LFG (Lexical Functional Grammar) as a syntactic formalism (output of a parse consists of c- and f-structures, no semantic representation yet)
- uses XLE (Xerox Linguistic Environment) as a parsing engine

- is aimed at developing a full-scale computational grammar for Georgian
- uses LFG (Lexical Functional Grammar) as a syntactic formalism (output of a parse consists of c- and f-structures, no semantic representation yet)
- uses XLE (Xerox Linguistic Environment) as a parsing engine
- uses visualization, treebanking and corpus tools developed at Aksis/University of Bergen

- is aimed at developing a full-scale computational grammar for Georgian
- uses LFG (Lexical Functional Grammar) as a syntactic formalism (output of a parse consists of c- and f-structures, no semantic representation yet)
- uses XLE (Xerox Linguistic Environment) as a parsing engine
- uses visualization, treebanking and corpus tools developed at Aksis/University of Bergen
- is part of the international Parallel Grammar (ParGram) project, which coordinates the development of LFG grammars in a parallel manner using XLE

Time frames and status quo:

started around 1995 with morphology, 2005 with syntax

#### Time frames and status quo:

- started around 1995 with morphology, 2005 with syntax
- most of the morphology is covered (missing: proper names; compounds)

#### Time frames and status quo:

- started around 1995 with morphology, 2005 with syntax
- most of the morphology is covered (missing: proper names; compounds)
- most basic syntactic constructions are covered (incomplete: subcategorization frames, constructions involving infinite forms, appositions, much more)

#### Time frames and status quo:

- started around 1995 with morphology, 2005 with syntax
- most of the morphology is covered (missing: proper names; compounds)
- most basic syntactic constructions are covered (incomplete: subcategorization frames, constructions involving infinite forms, appositions, much more)
- funding: none

### **Outline**

- 1 Introduction: The Georgian grammar project
- Lexical-Functional Grammar
- Morphology and Morphosyntax
- Some aspects of Georgian syntax
- 5 Tools for LFG grammar development

a generative linguistic framework



- a generative linguistic framework
- initiated by Joan Bresnan and Ronald Kaplan in the 1970s to overcome conceptual and explanatory shortcomings of Chomsky's transformational grammar

- a generative linguistic framework
- initiated by Joan Bresnan and Ronald Kaplan in the 1970s to overcome conceptual and explanatory shortcomings of Chomsky's transformational grammar
- constraint-based; no transformations

- a generative linguistic framework
- initiated by Joan Bresnan and Ronald Kaplan in the 1970s to overcome conceptual and explanatory shortcomings of Chomsky's transformational grammar
- constraint-based; no transformations
- rigid formalism, well-suited for implementation and efficient parsing, as well as for theoretical work

parallel description of linguistic entities by:

- parallel description of linguistic entities by:
  - C(onstituent)-structure: phrase structure tree (often modeled according to Xbar-syntax principles)

- parallel description of linguistic entities by:
  - C(onstituent)-structure: phrase structure tree (often modeled according to Xbar-syntax principles)
  - F(unctional)-structure: attribute-value matrix which recursively correlates the semantic argument structure of predicates with grammatical functions

- parallel description of linguistic entities by:
  - C(onstituent)-structure: phrase structure tree (often modeled according to Xbar-syntax principles)
  - F(unctional)-structure: attribute-value matrix which recursively correlates the semantic argument structure of predicates with grammatical functions

C- and f-structure are related by a projection relation



- parallel description of linguistic entities by:
  - C(onstituent)-structure: phrase structure tree (often modeled according to Xbar-syntax principles)
  - F(unctional)-structure: attribute-value matrix which recursively correlates the semantic argument structure of predicates with grammatical functions
  - C- and f-structure are related by a projection relation
- Structural relations are formally described by phrase structure rules annotated with functional equations

- parallel description of linguistic entities by:
  - C(onstituent)-structure: phrase structure tree (often modeled according to Xbar-syntax principles)
  - F(unctional)-structure: attribute-value matrix which recursively correlates the semantic argument structure of predicates with grammatical functions
  - C- and f-structure are related by a projection relation
- Structural relations are formally described by phrase structure rules annotated with functional equations
- Lexical items (lemmas and morphological features) are annotated with a lexical category and functional equations (including argument structures for verbs)

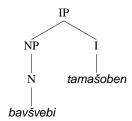
bavšvebi tamašoben. children they-play.

'The children are playing.'

bavšvebi tamašoben. children they-play.

'The children are playing.'

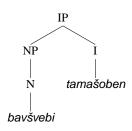
#### c-structure



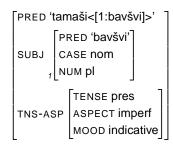
bavšvebi tamašoben. children they-play.

'The children are playing.'

#### c-structure

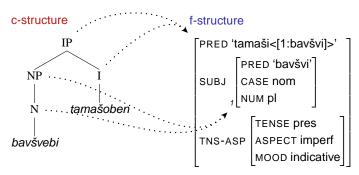


#### f-structure



bavšvebi tamašoben. children they-play.

'The children are playing.'





### Parsing steps:

Tokenization

- Tokenization
- Morphological analysis

- Tokenization
- Morphological analysis
- Lexical insertion

- Tokenization
- Morphological analysis
- Lexical insertion
  - lexemes and morphological features are entries in the LFG lexicon

- Tokenization
- Morphological analysis
- Lexical insertion
  - lexemes and morphological features are entries in the LFG lexicon
  - they are annotated with (sub)lexical categories and functional equations

- Tokenization
- Morphological analysis
- Lexical insertion
  - lexemes and morphological features are entries in the LFG lexicon
  - they are annotated with (sub)lexical categories and functional equations
  - these are used to initialize the parse chart

- Tokenization
- Morphological analysis
- Lexical insertion
  - lexemes and morphological features are entries in the LFG lexicon
  - they are annotated with (sub)lexical categories and functional equations
  - these are used to initialize the parse chart
- Chart parsing with phrase structure rules

- Tokenization
- Morphological analysis
- Lexical insertion
  - lexemes and morphological features are entries in the LFG lexicon
  - they are annotated with (sub)lexical categories and functional equations
  - these are used to initialize the parse chart
- Chart parsing with phrase structure rules
- Solving of functional equations

### **Outline**

- 1 Introduction: The Georgian grammar project
- Lexical-Functional Grammar
- Morphology and Morphosyntax
- Some aspects of Georgian syntax
- 5 Tools for LFG grammar development

### Parsing model

 First approach: Finite state transducer augmented with feature structure unification

- First approach: Finite state transducer augmented with feature structure unification
  - Disjunctive unification with a lexicon of existing forms to discard nonexisting verb analyses

- First approach: Finite state transducer augmented with feature structure unification
  - Disjunctive unification with a lexicon of existing forms to discard nonexisting verb analyses
  - Implemented in Common Lisp, based on Parc Xerox's old fsa module

- First approach: Finite state transducer augmented with feature structure unification
  - Disjunctive unification with a lexicon of existing forms to discard nonexisting verb analyses
  - Implemented in Common Lisp, based on Parc Xerox's old fsa module
- New implementation based on fst (Xerox finite state tool)

- First approach: Finite state transducer augmented with feature structure unification
  - Disjunctive unification with a lexicon of existing forms to discard nonexisting verb analyses
  - Implemented in Common Lisp, based on Parc Xerox's old fsa module
- New implementation based on fst (Xerox finite state tool)
  - automatically derived from old implementation

- First approach: Finite state transducer augmented with feature structure unification
  - Disjunctive unification with a lexicon of existing forms to discard nonexisting verb analyses
  - Implemented in Common Lisp, based on Parc Xerox's old fsa module
- New implementation based on fst (Xerox finite state tool)
  - automatically derived from old implementation
  - flag diacritics mimic feature structure unification; compiled out at the end ⇒ pure finite state

- First approach: Finite state transducer augmented with feature structure unification
  - Disjunctive unification with a lexicon of existing forms to discard nonexisting verb analyses
  - Implemented in Common Lisp, based on Parc Xerox's old fsa module
- New implementation based on fst (Xerox finite state tool)
  - automatically derived from old implementation
  - flag diacritics mimic feature structure unification; compiled out at the end ⇒ pure finite state
  - lexicon compiled into the transducer

- First approach: Finite state transducer augmented with feature structure unification
  - Disjunctive unification with a lexicon of existing forms to discard nonexisting verb analyses
  - Implemented in Common Lisp, based on Parc Xerox's old fsa module
- New implementation based on fst (Xerox finite state tool)
  - automatically derived from old implementation
  - flag diacritics mimic feature structure unification; compiled out at the end ⇒ pure finite state
  - lexicon compiled into the transducer
  - interfaces well with XLE



### The lexicon

 Derived from Kita Tschenkélis 'Georgisch-Deutsches Wörterbuch' (52 000 entries, 3 823 verb entries) and other sources (74 000 nouns and adjectives altogether)

### The lexicon

 Derived from Kita Tschenkélis 'Georgisch-Deutsches Wörterbuch' (52 000 entries, 3 823 verb entries) and other sources (74 000 nouns and adjectives altogether)

Typical analyses:

'wine'

 $\dot{g}vino \rightarrow \dot{g}vino+N+Nom+Sg$ 

```
Typical analyses:
```

```
'wine'
```

*ġvino* → ġvino+N+Nom+Sg

'for the girls, too'

gogo-eb-isa-tvis-ac → gogo+N+Anim+Full+Gen+Pl+Tvis+C

### Typical analyses:

```
'wine'
```

*ġvino* → ġvino+N+Nom+Sg

'for the girls, too'

gogo-eb-isa-tvis-ac → gogo+N+Anim+Full+Gen+Pl+Tvis+C

'in childhood'

bavšvob-isa-s → bavšvoba+N+DGen+DSg+Dat+Sg

```
Typical analyses:
'wine'
ġvino → ġvino+N+Nom+Sg
'for the girls, too'
gogo-eb-isa-tvis-ac → gogo+N+Anim+Full+Gen+Pl+Tvis+C
'in childhood'
bavšvob-isa-s → bavšvoba+N+DGen+DSg+Dat+Sg
'I apparently painted it'/'he will paint it for me'
da-mi-xat-av-s →
          { da-xatva-3569-5+V+Trans+Perf+Subj1Sg+Obj3
           da-xatva-3569-18+V+Trans+Perf+Subj1Sg+Obj3
           da-xatva-3569-18+V+Trans+Fut+Subj3Sq+Obj1Sq } s
```

Each verb lexeme in the LFG lexicon is associated with one or more subcategorization frames (argument structures) and a mapping of each of the arguments to a grammatical function (one of SUBJ(ect), OBJ(ect), OBJ(ect), OBJ(efficiary), OBL(ique), etc.).

Each verb lexeme in the LFG lexicon is associated with one or more subcategorization frames (argument structures) and a mapping of each of the arguments to a grammatical function (one of SUBJ(ect), OBJ(ect), OBJ(efficiary), OBL(ique), etc.).

ga-v-u-ket-eb 'I will do it for him/her':
 ga-keteba<agent, benefic, theme>

Each verb lexeme in the LFG lexicon is associated with one or more subcategorization frames (argument structures) and a mapping of each of the arguments to a grammatical function (one of SUBJ(ect), OBJ(ect), OBJ(ect), OBJ(efficiary), OBL(ique), etc.).

ga-v-u-ket-eb 'I will do it for him/her':

Each verb lexeme in the LFG lexicon is associated with one or more subcategorization frames (argument structures) and a mapping of each of the arguments to a grammatical function (one of SUBJ(ect), OBJ(ect), OBJ(efficiary), OBL(ique), etc.).

ga-v-u-ket-eb 'I will do it for him/her':

The verb classification in Tschenkéli's *Georgisch–deutsches Wörterbuch* could be used directly to automatically derive a preliminary version of the Georgian LFG verb lexicon

Each verb lexeme in the LFG lexicon is associated with one or more subcategorization frames (argument structures) and a mapping of each of the arguments to a grammatical function (one of SUBJ(ect), OBJ(ect), OBJ(efficiary), OBL(ique), etc.).

ga-v-u-ket-eb 'I will do it for him/her':

The verb classification in Tschenkéli's *Georgisch–deutsches Wörterbuch* could be used directly to automatically derive a preliminary version of the Georgian LFG verb lexicon

Example: Tschenkéli's class T<sup>3</sup> maps to the argument structure

P<SUBJ, OBJben, OBJ>



In many cases the correct frames are not (easily) deducible from Tschenkéli's classification:

verbs taking oblique or genitive arguments:

```
ča-tvla<SUBJ, OBJ, OBL<sub>adv</sub>> 'consider sb. to be sth.' 
še-šineba<SUBJ, OBJ<sub>gen</sub>> 'be afraid of sb./sth.'
```

- verbs taking oblique or genitive arguments:
   ča-tvla<SUBJ, OBJ, OBL<sub>adv</sub>> 'consider sb. to be sth.'
   še-šineba<SUBJ, OBJ<sub>gen</sub>> 'be afraid of sb./sth.'
- Class III verbs: can be transitive and intransitive

- verbs taking oblique or genitive arguments:
   ča-tvla<SUBJ, OBJ, OBL<sub>adv</sub>> 'consider sb. to be sth.'
   še-šineba<SUBJ, OBJ<sub>gen</sub>> 'be afraid of sb./sth.'
- Class III verbs: can be transitive and intransitive
- verbs taking clausal arguments

- verbs taking oblique or genitive arguments:
   ča-tvla<SUBJ, OBJ, OBL<sub>adv</sub>> 'consider sb. to be sth.'
   še-šineba<SUBJ, OBJ<sub>gen</sub>> 'be afraid of sb./sth.'
- Class III verbs: can be transitive and intransitive
- verbs taking clausal arguments
- morphological passives: can be passives or unaccusatives
   ga-ket-deba (mtavrobis mier): 'it will be done (by the government)'
   ga-keteba<OBL-AG, SUBJ>

```
da-brun-deba (*dedis mier): '(s)he will return' da-bruneba<SUBJ>
```

student-ma ceril-i mo-m-cer-a. student.ERG letter.NOM he-wrote-it-to\_me.

The student wrote me a letter.

mi-cera < SUBJ, OBJben, OBJ >

student-ma ceril-i mo-m-cer-a. student.ERG letter.NOM he-wrote-it-to\_me.

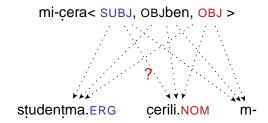
The student wrote me a letter.

mi-cera < SUBJ, OBJben, OBJ >

studentma.ERG cerili.NOM m-

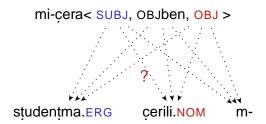
student-ma ceril-i mo-m-cer-a. student.ERG letter.NOM he-wrote-it-to\_me.

The student wrote me a letter.



student-ma ceril-i mo-m-cer-a. student.ERG letter.NOM he-wrote-it-to\_me.

The student wrote me a letter.



Georgian is head- and dependent-marking: verbal affixes and nominal case code grammatical functions.

## Case and affix alignment: Facts

Case alignment patterns

	SUBJ	OBJ	OBJben
Α	NOM	DAT	DAT
В	ERG	NOM	DAT
С	DAT	NOM	-tvis

## Case and affix alignment: Facts

Case alignment patterns

	SUBJ	OBJ	овJben
Α	NOM	DAT	DAT
В	ERG	NOM	DAT
С	DAT	NOM	-tvis

Alignment depending on verb class and tense group

	1	П	III	IV
	trans.	unacc.	unerg.	indir.
present	Α	Α	Α	С
aorist	В	Α	В	С
perfect	С	Α	С	С

## Case and affix alignment: Facts

Case alignment patterns

	SUBJ	OBJ	OBJben
Α	NOM	DAT	DAT
В	ERG	NOM	DAT
С	DAT	NOM	-tvis

Alignment depending on verb class and tense group

	ı	П	Ш	IV
	trans.	unacc.	unerg.	indir.
present	Α	Α	Α	С
aorist	В	Α	В	С
perfect	С	Α	С	С

Person/number affix alignment patterns

	SUBJ	OBJ	OBJben
A = B	v- (Fsubj)	<i>m</i> - (FовJ)	<i>h</i> - (FовJ)
С	<i>h-</i> (FовJ)	v- (Fsubj)	-

### Case and affix alignment: Implementation

Affixes: Alignment coded in the morphology

Example: 1st person plural (morphological) subject marker

ga-gv-i-ket-eb-i-a 'we apparently did it'

Functional equations are attached to morphology features:

```
+Subj1PI: (\uparrow SUBJ PERS) = 1

(\uparrow SUBJ NUM) = pI.

+Obj1PI: (\uparrow \_MORPH-SYNT \_AGR \_OBJ PERS) = 1

(\uparrow \_MORPH-SYNT \_AGR \_OBJ NUM) = pI.
```

### Case and affix alignment: Implementation

Case: Alignment coded in the syntax

Equations attached to verb lexicon entry

Example: transitive/unergative subject

### **Outline**

- 1 Introduction: The Georgian grammar project
- Lexical-Functional Grammar
- Morphology and Morphosyntax
- Some aspects of Georgian syntax
- 5 Tools for LFG grammar development

### Free word order

'Free word order' at the phrase level

#### 'Free word order' at the phrase level

 Subject and complements cannot be distinguished configurationally: no VP

#### 'Free word order' at the phrase level

- Subject and complements cannot be distinguished configurationally: no VP
- Finite verb and other constituents can occur in almost arbitrary order; or: an arbitrary permutation of the toplevel constituents of a grammatical sentence results in a grammatical sentence with the same propositional truth value

#### 'Free word order' at the phrase level

- Subject and complements cannot be distinguished configurationally: no VP
- Finite verb and other constituents can occur in almost arbitrary order; or: an arbitrary permutation of the toplevel constituents of a grammatical sentence results in a grammatical sentence with the same propositional truth value
- This is to be expected: since grammatical functions are coded morphologically, there is no need to repeat the coding configurationally

#### 'Free word order' at the phrase level

- Subject and complements cannot be distinguished configurationally: no VP
- Finite verb and other constituents can occur in almost arbitrary order; or: an arbitrary permutation of the toplevel constituents of a grammatical sentence results in a grammatical sentence with the same propositional truth value
- This is to be expected: since grammatical functions are coded morphologically, there is no need to repeat the coding configurationally
- ⇒ First approximation:

 $S \rightarrow V, XP^*$ 



Position is significant for Information structure: it is used to code the discurse functions FOCUS and TOPIC.

• FOCUS: immediately in front of inflected verb or in last position

- FOCUS: immediately in front of inflected verb or in last position
- TOPIC: initial position(s) to the left of FOCUS and verb

- FOCUS: immediately in front of inflected verb or in last position
- TOPIC: initial position(s) to the left of FOCUS and verb
- ⇒ Revision of phrase structure rules (compliant with Xbar theory, Bresnan 2001):

- FOCUS: immediately in front of inflected verb or in last position
- TOPIC: initial position(s) to the left of FOCUS and verb
- $\Rightarrow$  Revision of phrase structure rules (compliant with Xbar theory, Bresnan 2001):

I 
$$ightarrow$$
 finite verb IP  $ightarrow$  (XP) I' IP  $ightarrow$  XP IP  $ightarrow$  XP+

I is the category of the finite (inflected) verb:

I → finite verb



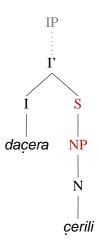
The nonprojective category S is the Complement of I:

$$I' \rightarrow I(S)$$

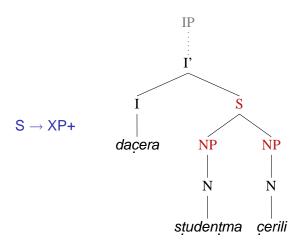


S contains all material to the right of the verb:

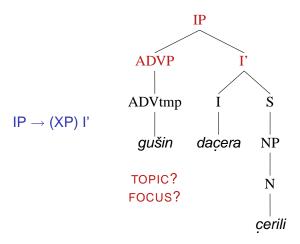




S contains all material to the right of the verb:

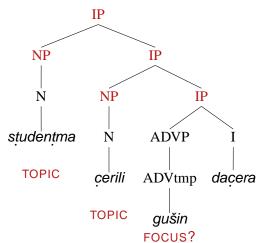


The Specifier of I is often TOPIC or FOCUS position:

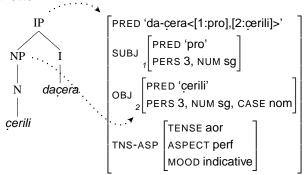


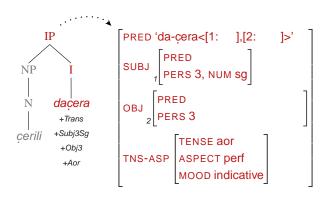
#### Material adjoint to IP is TOPIC:

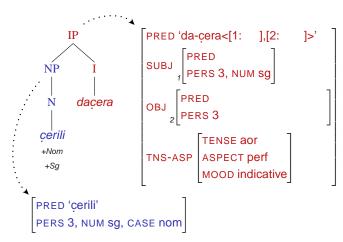


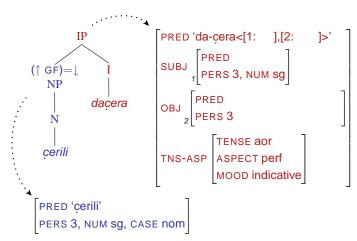


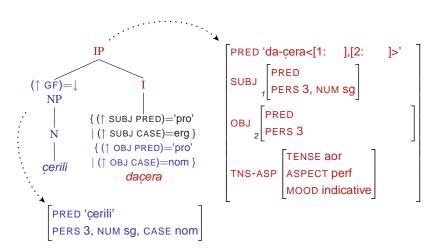
cerili dacera.
letter.NOM he-wrote-it.AOR

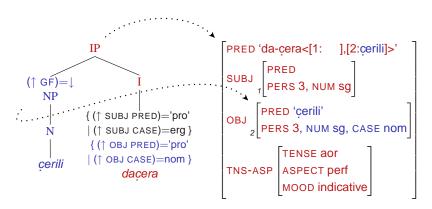


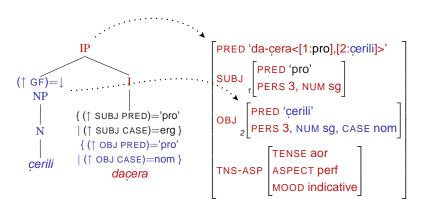






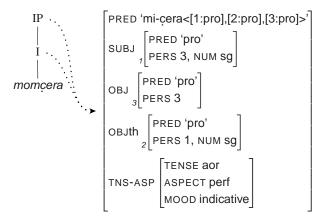






## Pro-drop: Example

'She wrote it to me.'



Predicate coordination: two verbs coordinate at V level (as opposed to sentence coordination with ellipsis), and share their arguments.

Predicate coordination: two verbs coordinate at V level (as opposed to sentence coordination with ellipsis), and share their arguments.

Crosslinguistically common restriction: Both verbs should assign the same semantic roles and grammatical functions to their shared arguments.

Predicate coordination: two verbs coordinate at V level (as opposed to sentence coordination with ellipsis), and share their arguments.

- Crosslinguistically common restriction: Both verbs should assign the same semantic roles and grammatical functions to their shared arguments.
- Frequent additional restriction: The case of a common argument should be licensed by both verbs, or, when there is case syncretism, be compatible with both verbs' requirements.

Predicate coordination: two verbs coordinate at V level (as opposed to sentence coordination with ellipsis), and share their arguments.

- Crosslinguistically common restriction: Both verbs should assign the same semantic roles and grammatical functions to their shared arguments.
- Frequent additional restriction: The case of a common argument should be licensed by both verbs, or, when there is case syncretism, be compatible with both verbs' requirements.

In Georgian: less restrictive conditions:

Predicate coordination: two verbs coordinate at V level (as opposed to sentence coordination with ellipsis), and share their arguments.

- Crosslinguistically common restriction: Both verbs should assign the same semantic roles and grammatical functions to their shared arguments.
- Frequent additional restriction: The case of a common argument should be licensed by both verbs, or, when there is case syncretism, be compatible with both verbs' requirements.

#### In Georgian: less restrictive conditions:

The case of an argument needs only to be licensed by the verb nearest to it.

mas [uqvars da apasebs] tavis meuġle-s. he.DAT loves and esteems his-own wife.DAT

'He loves and esteems his own wife.'

mas [uqvars da apasebs] tavis meuġle-s. he.DAT loves and esteems his-own wife.DAT

'He loves and esteems his own wife.'

	uqvars (IV)	apasebs (I)
	'he loves her'	'he esteems her'
thematic roles	< exp, theme >	< exp, theme >
functions	SUBJ, OBJ	SUBJ, OBJ
case marking	DAT, NOM	NOM, DAT

```
*me [miqvars da mapasebs] čemi meuġle.

I I-love-her and she-esteems-me my wife.NOM
```

'I [love, and am esteemed by,] my own wife.'

```
*me [miqvars da mapasebs] čemi meuġle.

I I-love-her and she-esteems-me my wife.NOM
```

'I [love, and am esteemed by,] my own wife.'

	uqvars (IV) 'he loves her'	apasebs (I) 'he esteems her'
thematic roles	< exp, theme >	< exp, theme >
functions	SUBJ, OBJ	SUBJ, OBJ
case marking	DAT, NOM	NOM, DAT
	· • • • • • • • • • • • • • • • • • • •	*

??misi moġvaçeoba mosçons da pasdeba xalxis mier. his public-activity it-likes-it and it-is-esteemed the-people by

'His public activity is liked and esteemed by the people.'

??misi moġvaçeoba mosçons da pasdeba xalxis mier. his public-activity it-likes-it and it-is-esteemed the-people by

'His public activity is liked and esteemed by the people.'

	moscons (IV) 'he likes it'	pasdeba (II) 'it is esteemed'
case marking	DAT, NOM	mier, NOM
functions	SUBJ, OBJ	OBL, SUBJ
thematic roles	< exp, theme >	< exp, theme >
	· • •	
		*

## Predicate coordination: Implementation

Simple case: both verbs license the same cases



## Predicate coordination: Implementation

Simple case: both verbs license the same cases

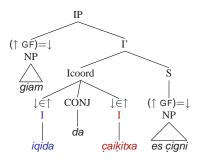
giam iqida da çaikitxa es çigni. Gia bought and read this book.

Simple case: both verbs license the same cases

giam iqida da çaikitxa es çigni. Gia bought and read this book.

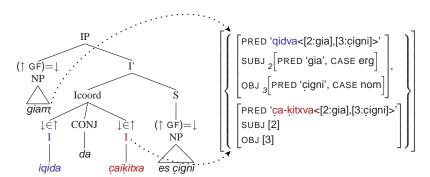
#### Simple case: both verbs license the same cases

giam iqida da çaikitxa es çigni. Gia bought and read this book.



#### Simple case: both verbs license the same cases

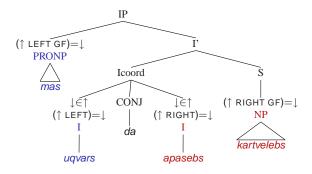
giam iqida da çaikitxa es çigni. Gia bought and read this book.



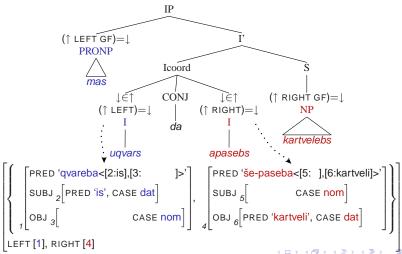
General case: the verbs license different cases

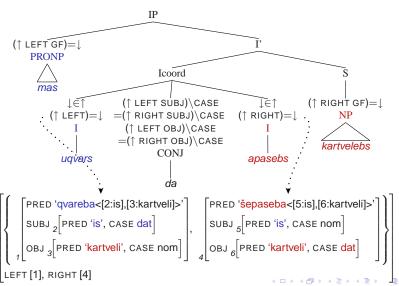


#### General case: the verbs license different cases



#### General case: the verbs license different cases

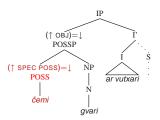


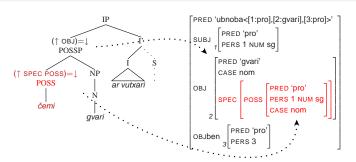


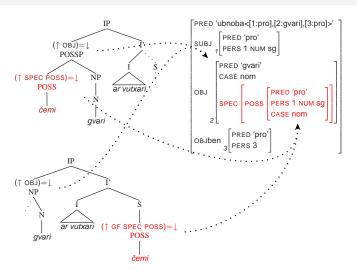
```
čem-i gvar-i ar v-u-txar-i [[my.NOM].POSS last-name.NOM].NP not I.to-him.told.it.
```

'I did not tell him my last name.'

```
gvar-i ar v-u-txar-i čem-i [last-name.NOM].NP not I.to-him.told.it [my.NOM].POSS.
```







## **Outline**

- 1 Introduction: The Georgian grammar project
- Lexical-Functional Grammar
- Morphology and Morphosyntax
- Some aspects of Georgian syntax
- Tools for LFG grammar development

### XLE (Xerox Linguistic Environment)

 is a sophisticated development platform for LFG grammars developed by the Palo Alto Research Center with active participation of some of the inventors of LFG.

- is a sophisticated development platform for LFG grammars developed by the Palo Alto Research Center with active participation of some of the inventors of LFG.
- consists of a parser, a generator and a transfer module

- is a sophisticated development platform for LFG grammars developed by the Palo Alto Research Center with active participation of some of the inventors of LFG.
- consists of a parser, a generator and a transfer module
- can be used both from Emacs via a Tcl/Tk interface that provides powerful viewing and debugging facilities, and as a shared library, which opens up for integrating XLE into custom software

- is a sophisticated development platform for LFG grammars developed by the Palo Alto Research Center with active participation of some of the inventors of LFG.
- consists of a parser, a generator and a transfer module
- can be used both from Emacs via a Tcl/Tk interface that provides powerful viewing and debugging facilities, and as a shared library, which opens up for integrating XLE into custom software
- Tokenization and morphological analysis is normally done with the Xerox finite state tool, fst

#### XLE-Web

 easy-to-use pedagogical Web interface to XLE for parsing sentences on the fly

- easy-to-use pedagogical Web interface to XLE for parsing sentences on the fly
- in use for several of the ParGram grammars (among others Norwegian, English, German, Welsh and Malagassy)

- easy-to-use pedagogical Web interface to XLE for parsing sentences on the fly
- in use for several of the ParGram grammars (among others Norwegian, English, German, Welsh and Malagassy)
- display of c- and f-structures of LFG analyses

- easy-to-use pedagogical Web interface to XLE for parsing sentences on the fly
- in use for several of the ParGram grammars (among others Norwegian, English, German, Welsh and Malagassy)
- display of c- and f-structures of LFG analyses
- visualization of the mapping from c- to f-structure

- easy-to-use pedagogical Web interface to XLE for parsing sentences on the fly
- in use for several of the ParGram grammars (among others Norwegian, English, German, Welsh and Malagassy)
- display of c- and f-structures of LFG analyses
- visualization of the mapping from c- to f-structure
- display of compact packed representations of c- and f-structures that combine the c- resp. f-structures of all analyses of a given parse into one c- resp. f-structure graph

Tasks when developing a large grammar:

In order to monitor progress, to assess coverage and to compare analyses across different grammar versions:

Tasks when developing a large grammar:

In order to monitor progress, to assess coverage and to compare analyses across different grammar versions:

run the grammar on a set of sample sentences

Tasks when developing a large grammar:

In order to monitor progress, to assess coverage and to compare analyses across different grammar versions:

- run the grammar on a set of sample sentences
- store the parse results

Tasks when developing a large grammar:

In order to monitor progress, to assess coverage and to compare analyses across different grammar versions:

- run the grammar on a set of sample sentences
- store the parse results
- rerun successive versions of the grammar on the same sentences and compare the results

When the grammar has reached acceptable coverage, one wants to:

When the grammar has reached acceptable coverage, one wants to:

 run the grammar on a larger set of sentences (perhaps chosen from running text)

When the grammar has reached acceptable coverage, one wants to:

- run the grammar on a larger set of sentences (perhaps chosen from running text)
- develop a treebank in the sense of a linguistic resource

When the grammar has reached acceptable coverage, one wants to:

- run the grammar on a larger set of sentences (perhaps chosen from running text)
- develop a treebank in the sense of a linguistic resource

Problems:

When the grammar has reached acceptable coverage, one wants to:

- run the grammar on a larger set of sentences (perhaps chosen from running text)
- develop a treebank in the sense of a linguistic resource

#### Problems:

 sentences of only moderate complexity often are highly ambiguous

When the grammar has reached acceptable coverage, one wants to:

- run the grammar on a larger set of sentences (perhaps chosen from running text)
- develop a treebank in the sense of a linguistic resource

#### Problems:

- sentences of only moderate complexity often are highly ambiguous
- the desired or correct reading is only one of the analyses offered by the grammar

When the grammar has reached acceptable coverage, one wants to:

- run the grammar on a larger set of sentences (perhaps chosen from running text)
- develop a treebank in the sense of a linguistic resource

#### Problems:

- sentences of only moderate complexity often are highly ambiguous
- the desired or correct reading is only one of the analyses offered by the grammar
- ⇒ Need for manual disambiguation of the parses in an efficient way

### LFG Parsebanker

LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

### LFG Parsebanker

LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

 developed in the Trepil project (together with Victoria Rosén and Koenraad de Smedt, Bergen)

LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

- developed in the Trepil project (together with Victoria Rosén and Koenraad de Smedt, Bergen)
- originally for Norwegian, but language independent

LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

- developed in the Trepil project (together with Victoria Rosén and Koenraad de Smedt, Bergen)
- originally for Norwegian, but language independent

# LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

- developed in the Trepil project (together with Victoria Rosén and Koenraad de Smedt, Bergen)
- originally for Norwegian, but language independent

Supports a process flow involving

automatic parsing with XLE

# LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

- developed in the Trepil project (together with Victoria Rosén and Koenraad de Smedt, Bergen)
- originally for Norwegian, but language independent

- automatic parsing with XLE
- viewing with XLE-Web

# LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

- developed in the Trepil project (together with Victoria Rosén and Koenraad de Smedt, Bergen)
- originally for Norwegian, but language independent

- automatic parsing with XLE
- viewing with XLE-Web
- structural c- and f-structure queries based on the TIGERSearch treebank search tool

# LFG Parsebanker: Web-based toolkit for building and manual disambiguation of an LFG treebank

- developed in the Trepil project (together with Victoria Rosén and Koenraad de Smedt, Bergen)
- originally for Norwegian, but language independent

- automatic parsing with XLE
- viewing with XLE-Web
- structural c- and f-structure queries based on the TIGERSearch treebank search tool
- efficient manual disambiguation by means of discriminants

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

Our discriminants:

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

#### Our discriminants:

specifically designed for LFG grammars

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

- specifically designed for LFG grammars
- four major types: lexical, morphological, c-structure and f-structure discriminants

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

- specifically designed for LFG grammars
- four major types: lexical, morphological, c-structure and f-structure discriminants
- A lexical discriminant is a word form together with its part of speech

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

- specifically designed for LFG grammars
- four major types: lexical, morphological, c-structure and f-structure discriminants
- A lexical discriminant is a word form together with its part of speech
- A morphological discriminant is a base form with the tags it receives from morphological preprocessing

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

- specifically designed for LFG grammars
- four major types: lexical, morphological, c-structure and f-structure discriminants
- A lexical discriminant is a word form together with its part of speech
- A morphological discriminant is a base form with the tags it receives from morphological preprocessing
- C-structure discriminants are based on minimal subtrees, a minimal subtree being defined as a mother node and her daughters

Discriminant: 'Any elementary linguistic property of an analysis that is not shared by all analyses' (David Carter).

- specifically designed for LFG grammars
- four major types: lexical, morphological, c-structure and f-structure discriminants
- A lexical discriminant is a word form together with its part of speech
- A morphological discriminant is a base form with the tags it receives from morphological preprocessing
- C-structure discriminants are based on minimal subtrees, a minimal subtree being defined as a mother node and her daughters
- F-structure discriminants are based on partial paths through f-structures

An indispensable resource for research in Georgian syntax is a searchable text corpus of decent size.

Available text collections on the Internet:

• the electronic newspaper archive Opentext (> 75 million words)

An indispensable resource for research in Georgian syntax is a searchable text corpus of decent size.

Available text collections on the Internet:

- the electronic newspaper archive Opentext (> 75 million words)
- the text archive of Radio tavisupleba (the Georgian service of Radio Free Europe/Radio Liberty) (around eight million words)

An indispensable resource for research in Georgian syntax is a searchable text corpus of decent size.

Available text collections on the Internet:

- the electronic newspaper archive Opentext (> 75 million words)
- the text archive of Radio tavisupleba (the Georgian service of Radio Free Europe/Radio Liberty) (around eight million words)
- fictional texts (both prose and poetry): the UNESCO project Digital collection of Georgian classical literature (three million words)

An indispensable resource for research in Georgian syntax is a searchable text corpus of decent size.

Available text collections on the Internet:

- the electronic newspaper archive Opentext (> 75 million words)
- the text archive of Radio tavisupleba (the Georgian service of Radio Free Europe/Radio Liberty) (around eight million words)
- fictional texts (both prose and poetry): the UNESCO project Digital collection of Georgian classical literature (three million words)

I have harvested these text collections and imported them into corpus query software based on Corpus Workbench (IMS Stuttgart) which is being developed at Aksis.