

Please cite as:

Rogier Blokland, Marina Fedina, Niko Partanen, Michael Rießler. 2019. "Closing the cycle of language documentation, corpus building and corpus-based description of Zyrian Komi." Paper held at the conference *Descriptive grammars and typology. The challenges of writing grammars of underdescribed and endangered languages, University of Helsinki, 27–29 March 2019.* Date: 28.03.2019. (Version of 17.04.2019). Licence: CC-BY-ND





Kotimaisten kielten keskus INSTITUTET FÖR DE INHEMSKA SPRÅKEN

JTE FOR THE LANGUAGES OF FINLAND



Closing the cycle of language documentation, corpus building and corpus-based description of Zyrian Komi



 Rogier Blokland (Uppsala University, Sweden)
 Marina Fedina (The Finno-Ugric Laboratory for Support of the Electronic Representation of Regional Languages, Syktyvkar, Komi Republic, Russia)
 Niko Partanen (Institute for the Languages of Finland, Helsinki, Finland)
 Michael Rießler (Bielefeld University, Germany)





Our Komi projects

1st project ('Izva Komi: Building an Annotated Digital Corpus for Future Research on Komi Speech Communities in Northernmost Russia'; 2014–2016)

- **Fieldwork** in several locations, inside and outside of the Komi Republic
- ~45 hours **transcribed**, 100 speakers (~80 hours recorded by us)
- ~400,000 tokens aligned at utterance level

2nd project ('Language Documentation meets Language Technology: The Next Step in the Description of Komi'; 2017–)

- Transcribed corpus > to Kielipankki
- Annotated portion in **Universal Dependencies treebank**
- Use of language technology instead of manual annotation



Our Komi grammar project

- Descriptive grammar with a **focus on syntax**
- Will be available **online**, written under **version control**
- All glossed examples in the grammar **connected to corpora**
- Allows examining occurrences of a specific phenomenon in corpora
- In addition to writing a 'traditional' grammar, another goal of ours is to **implement the grammar** with language technology





Background Zyrian Komi, our test language



Source: Komi mu newspaper, 6.12.2018

Source: OpenStreetMap.org

- Uralic language, spoken by about 160 000 speakers in northern Russia
- Relatively vital, but definitely endangered due to language shift to Russian
- Existing linguistic descriptions focus on phonology and morphology
- Orthographic standard, in which much material has been published
- Interest in preservation and further development, also in creating **language technology** for teaching and to **increase its functionality in the digital age**



Background NLP and Language Technology

Computational linguistics (NLP) = the analysis/generation of natural language

Language technology = the practical application of NLP

- spell checkers
- grammar taggers
- machine translation
- etc.



(North Saami; http://divvun.no/korrektur/speller-demo.html)





Workflow

- Our work relies on the NLP infrastructure available for Northern Eurasian endangered languages at **Giellatekno**. At the beginning of our workflow are **speech recordings with aligned transcriptions in ELAN**
- Finite-State morphology (FST) is used for rule-based modeling of stems and segmental affixes, as well as complex morphophonological rules
- Additional rules following **Constraint Grammar (CG)** are implemented for syntactic disambiguation and tagging dependency relations
- The **source code and documentation** is being developed using an SVN versioning system and is available under a GNU public license



Workflow: illustration

- Fieldwork (or "archive work", digitization, etc.)
- ELAN, transcription, translation
- Writing rules FST, CG, maintaining lexicon
- Applying grammar into texts in ELAN
 - Manual corrections for selected texts
 - Improving the grammar





Workflow: example

•									뛽 ELAN	5.3 - kpv_izva2	0140404lgu	evJA-fragment.eaf					
F	ile E	dit Annota	tion Tier	Туре	Search	View	Options	Window	Help								
										Grid	Text Su	btitles Lexicon	Comments	Recognizer	rs Metadata 🤇	Controls	
2								Volume: 100 kp	0 v_izva20144 Mute S v_izva20144 Mute S	04041gusevJA- iolo 04041gusevJA. iolo	-fragment.wa m4v		' <mark> </mark> 50 ' 25 ' 25	1 I I	1 1 1 50 1 1 1 50 50	1 1 1 75 1 1 75 1 1	100 100 100
	$\begin{array}{c c c c c c c c c c c c c c c c c c c $																
500								LINE M.L. LILES									
C	hpv_jzva2																
•	- part	0:5	0 00		00:	:00:55.50	00		00:00:56.	000		00:00:56.500		00:00:57.00	0	00:00:57.500	
6		IALM-1939	kpv_izva2	20140404lgu	usevJA-b-38	31											
	(21)	n@JAI-M-1	Me		пö			ciec		CÖBCEM	. ər		любит	12		висьтало	
-	12 [2	ord@JAI-M	ме		пö			ciec		совсем	03		любитны			висьтавны	
-	6	pos@JAI	Pron		CON	J	Pcle	?		Adv	v		v	CL	.В	V	CLB
		morph	Pers+Sg1	+N Pers+S	g1+C			_			Ne	g+Ind+Prt Neg+Ind+Pr	t TV+ConNeg T	TV+Imprt+S		TV+Ind+Prs+Sg3	
Image: State																	
d		NTP-M-198															



FST / CG

We apply grammar-based ("symbolic") NLP: the linguist writes a **formalized machine-readable** version of the grammar, and compiles it into a program capable of analyzing (and also generating) text input.

Finite-state transducer: modeling stems, affixes, linear morphology (Shoebox/Toolbox/FLEx/ELAN do this, too)



Constraint Grammar: disambiguation and dependency tagging (Shoebox/Toolbox/FLEx/ELAN cannot do this!)



FST / CG: Me yepu or cëŭ. 'I don't eat fish.'

Tagging ("glossing")

сёй <mark>сёй</mark>+N+Sg+Nom

сёй <mark>сёйны</mark>+V+ConNeg

сёй <mark>сёйны</mark>+V+Imprt+Sg2



FST / CG: Me yepu or cëŭ. 'I don't eat fish.'

Tagging ("glossing")

сёй <mark>сёй</mark>+N+Sg+Nom сёй <mark>сёйны</mark>+V+ConNeg

сёй <mark>сёйны</mark>+V+Imprt+Sg2

Disambiguation e.g. "IFF Rule": ConNeg if Neg to the left

сёй сёй+N+Sg+Nom сёй <mark>сёйны</mark>+V+ConNeg сёй <mark>сёйны</mark>+V+Imprt+Sg2



FST / CG: Me yepu or cëŭ. 'I don't eat fish.'

Full analysis (incl. dependency structure)

ме	ме	+Pron+Pers+Sg1+Nom	@SUBJ>	#1->3
чери	чери	+N+Sg+Nom	@OBJ>	#2->4
OF	03	+V+Neg+Ind+Prs+Sg1	@FAUX	#3->0
сёй	сёйны	+V+ConNeg	0 IMV	#4->3
•	•	+CLB		#5->3





[кру] Тайö нигаыс сетас сöмын ичöтик юкöн ывлавыв велöдысьлы.

'This book imparts only some of the knowledge needed by the student of the outdoors'



"<Тайö>" "тайо" Pron Dem Sg Nom @X #1->2 "<нигаыс>" "нига" N Sg Nom PxSg3 @SUBJ #2->3 "<сетас>" "сетны" V TV Ind Fut Sg3 @X #3->0 "<сомын>" "сомын" Adv @X #4->0 "<ичöтик>" "ичöтик" A Sg Nom @A< #5->6 "<юкöн>" "юкöн" N Sg Nom @X #6->0 "<ывлавыв>" "ывлавыв" Der Der/выв N Sg Nom @N< #7->8 "<велодысьлы>" "велöдысь" N Sg Dat @N< #8->8 "< >" "." CLB #9->9



Buryat Cantonese Catalan Chinese · Any **Classical Chinese** 衋 Coptic -Croatian Czech Danish Dutch English Erzya Estonian Faroese Finnish French X Galician German 讄 Gothic #= Greek 0 Hebrew . Hindi **Hindi English** . Hungarian Indonesian Irish Italian ٠ Japanese Kazakh Komi Zyrian :•: Korean

0

Kurmanji

home edit page issue tracker

This page pertains to UD version 2.

Universal Dependencies

ADV: adverb

Adverbs are words that typically modify adjectives, verbs or other adverbs for such categories as time, place, direction or manner.

Examples

- [kpv] *6vpa* "well"
- [kpv] дзик "completely; really"
- [kpv] vна "much: many: a lot"

ADV in other languages: [bg] [bm] [ca] [cs] [da] [en] [es] [et] [eu] [fi] [fro] [fr] [ga] [grc] [hu] [hy] [it] [ja] [kk] [myy] [no-overview] [pcm] [pt] [ru] [sl] [sv] [tr] [uk] [u] [urj] [yue] [zh]

Treebank Statistics: UD_Komi_Zyrian-Lattice: POS Tags: ADV

There are 96 ADV lemmas (12%), 99 ADV types (9%) and 202 ADV tokens (10%). Out of 15 observed tags, the rank of ADV is: 4 in number of lemmas, 4 in number of types and 4 in number of tokens.

The 10 most frequent ADV lemmas: нин, зэв, на, сомын, кыдзи, пыр, сідз, весиг, одйо, оні

The 10 most frequent ADV types: нин, зэв, на, сомын, кыдзи, пыр, сідз, весиг, оні, кыдз

The 10 most frequent ambiguous lemmas: сомын (ADV 7, PART 1), пыр (ADV 6, ADP 1), водз (ADP 3, ADV 3), друг (ADV 3, NOUN 1), медся (ADV 2, PART 2), 6öpuh (ADP 1, ADV 1), KOP (SCONJ 3, ADV 1, NOUN 1), KYTIIOM (ADV 1, PRON 1), KÖTE (ADV 1, PART 1, SCONJ 1), MEA (SCONJ 2, ADV 1)

The 10 most frequent ambiguous types: сомын (ADV 7, PART 1), друг (ADV 2, NOUN 1), медся (ADV 2, PART 2), Мыйла (ADV 1, SCONJ 1), борын (ADP 1, ADY 1), KOP (SCONJ 3, ADY 1), KOTH (ADV 1, SCONJ 1), MCJ (ADV 1, SCONJ 1), MO3 (ADP 1, ADV 1), CTAB (DET 4, PRON 3, ADV 1)

- сёмын
 - ADV 7: Сэк жö удиті на воны сомын пипу рас весьтодз.
 - <u>ракт</u> 1: Медым содтыны урожай идраломын одъяс, коло не сомын сувтодны удж выло став вундан машинаяс, но и приспособитны лобогрейкаё ытшкан машинаяс да используйтны найёс тыр нагрузкаён урожай идралём вылё.

друг

- ADV 2: Модлаполысь Присада ді вывса пожома яг весьтын друг тыдовтчис эшкын кодь кымор пласт, коді син водзын кутіс быдмыны.
- NOUN 1: Да мый тэ, друг, мый тэ!



Explicit references to sources, versions & licenses

	🗶 ELAN 4	9.3 - kpv_izva20150703-01-b.eaf	
		Grid Text Subtitles Lexicon Comments Recognizers	Controls
		Select Metadata Source No metadata source selected.	Configure
	0:04:54.375 Selector: 00:04:54.375	-000456510 2135 ← → ↓ ↑ Selection Mode Loop Mode ↓0	A ALARIA DUMERANIMATINA CONTRACTOR
kpv_izva20150703-01-b.wav 0	00:04:55.000 00:04:56.000	00:04:57.000 00:04:58.000 00:04:59.000 00:05:00.000	00:05:01.000
	Unsupported compression type for Wave file		
*			
ft-eng@MSF-F-1968 [1]	00:04:55.000 00:04:56.000	00:04:57.000 00:04:58.000 00:04:59.000 00:05:00.000	00:05:01.000
ft-rus@MSF-F-1968 (1)			
ref@VPC-M-1993 [51]	kpv_izva20150703-01-b-099		kpv_izva20150703-01-b-10
orth@VPC-M-1993 [51]	Кутшемке сьёкыдлунъяс тэнад вёлісны из?		А тэ наес понимайтін?
@- word@VPC-M-1993 [309]	Кутшем сьокыд тэнад волісны из ?		А тэ наес пони ?
- ft-eng@VPC-M-1993 (51)	So did you have some kind of difficulties?		But you understood them?
- fl-rus@VPC-M-1993 [51]	Какие-то трудности у тебя были нет?		А ты их понимала?
	in a second s	kpv izva20150703-01-b-100	kpv_izva20
ref@IGT-F-1996 [93]	150703-		
ref@IGT-F-1996 [93] P_ orth@IGT-F-1996 [93]	<u>J150703-</u>	Сьокыдлун то, что наа простэ менэ из понимайтныс водздзык.	Ме ны пон
ref@IGT-F-1996 [93] orth@IGT-F-1996 [93] - word@IGT-F-1996 [806]	1150703-	Събкыдлун то, что наа прёстэ менэ из пёнимайтныс водадзык. Събкыд то , что наа прёстэ менэ из пёнима водада ,	Ме ны пон Ме ны
 ref@IGT-F-1996 [83] orth@IGT-F-1996 [83] word@IGT-F-1996 [806] f-eng@IGT-F-1996 [83] 	7150703-	Сьбжыдлун то, что наа простэ менэ из понимайтные водздзык. Събжыд то , что наа простэ менэ из понима, водздз , Difficulties with that they just didn't understand me earlier.	Ме ны пон Ме ны

Source: IKDP corpus, the Language Bank version 1.0 (coming soon) kpv_izva20150703-01-b Чељад школајасын сомын вермасны зев ічотіка то́дмавны, сені оз удітны велавны вело́дчыејас—кыз ескої колої то́дмавны возої. Сыко́д поъої то̀дмаены со́мын гырые школајасын—Універеітетјасын. Сіјої школајасас емоїс аслыс полоїс јо́рто̀дјас-керкајас. Сіјої jо́рто̀дјасас чуко́рто́мао́с быdеама полоїс ужалан ко́лујјассої. Сіјої jо́рто̀дјасыс шусо́ны л абор ат ор і ј а ј а с ої н. Быdмоїгјас, немоїсјас вело́дны емоїс во́рјасын зорізјас; jyjac, тыјас-бердын аслыс полоїс керкајас (лабораторіјајас, стантсіјајас).

Почо зев унатор ывлавылые то́дмавны лабораторіјајасын, ещанік колујон—сы-понда ковмас то̀дмаены војдор нігако́д. И іга отсалас го́го́рвоны мед-гырыс ывлавыв петко́дчо́м јассо. Ніга індас кор, кыті ывлавылые почо то̀дмавны. Січіко́н почо́ зев уна то̀дмавны лабораторіјајасын. Ставсо дерт он понды то̀дны. Тајо нігаыс гіжо́ма сещо́м јо̀злы. Унао́н оз вермыны вескавны некущо́м лабораторіјао́—налы отсо́г-пыд-di лоо́ тајо ніга.

Тані абу став пічіка, кіміја, біологіја. Тані емёс сёмын мед-колан торјас. Тајё нігаыс сетас сёмын ічэтік јукён ывлавыв велёдыслы.

Source: Fenno-Ugrica collection http://urn.fi/URN:NBN:fi-fe2014070132058



Pros and cons

	Sketch grammar	Implemented grammar
Learning curve	V Common methodology	X Steep learning curve (command line tools)
Application	Cannot be applied as a corpus annotation tool	Can be directly applied as a corpus annotation tool
Time investment	Describing rules and paradigms is time-consuming	Programming grammar-based NLP is time-consuming
Audience of the grammar	Comparative linguists, specialists in the respective language (perhaps also community members, teachers, etc.)	Computational linguists (less likely to be community members, teachers, etc.)



Summary

- Our approach is fundamentally different **in practice** from common grammaticographical approaches in Documentary Linguistics
- It is not different in its aims:
 - ensuring that there are resources for this language (basic documentation)
 - providing sufficient analyses for linguistic research
 - creating something that is useful for the community (> creates preliminaries for digital infrastructure)
- In its **results**, it goes beyond common approaches and **links Documentary Linguistics closer** to established NLP methods in corpus linguistics



References

- Blokland, Rogier; Niko Partanen, Michael Rießler, & Joshua Wilbur. 2019. Using computational approaches to integrate endangered language legacy data into documentation corpora. *Proceedings of ComputEL-3 2019*, 24–30. <u>https://scholar.colorado.edu/scil-cmel/vol2/iss1/5</u>
- Gerstenberger, Ciprian; Niko Partanen, Michael Rießler, & Joshua Wilbur. 2016. Utilizing language technology in the documentation of endangered Uralic languages. *Northern European Journal of Language Technology*, 29–47. <u>http://www.nejlt.ep.liu.se/2016/v4/a03</u>
- Beesley, Kenneth & Lauri Karttunen. 2003. Finite State Morphology. CSLI.
- Fedina, Marina & Olga Kuzivanova. 2013. Le statut officiel et social du komi sur le territoire de la République de Komi. Histoire et situation actuelle. *Études finno-ougriennes* 45, 1–13.
- Karlsson, Fred; Arto Voutilainen, Juha Heikkilä, & Arto Anttila (eds). 1995. Constraint Grammar: A language-independent system for parsing unrestricted text. Mouton de Gruyter.
- Partanen, Niko; Rogier Blokland, KyungTae Lim, Thierry Poibeau, & Michael Rießler. 2018. The first Komi-Zyrian Universal Dependencies treebanks. *Proceedings of UDW 2018*, 126–132. <u>https://aclanthology.info/papers/W18-6015/w18-6015</u>
- Partanen, Niko; KyungTae Lim, Michael Rießler, & Thierry Poibeau. 2018. Dependency parsing of code-switching data with cross-lingual feature representations. *Proceedings of IWCLUL 2018*, 1–17. https://aclweb.org/anthology/papers/W/W18/W18-0201
- Trosterud, Trond. 2006. Grammatically based language technology for minority languages. *Lesser-known languages of South Asia*, ed. by Anju Saxena et al. Mouton de Gruyter, 293–316.





Аттьö! Kiitos! Thank you!



