

Zsófia Schön

*OUIDB Yugan Khanty corpus*



Second workshop on Uralic prosody

RIL HAS

Budapest, 28.09.2017



# Outline

1. Project Background
2. Fieldwork Recordings
3. Corpus

# 1. Project Background



# 1.1 Project Background



## OUDB

- Ob-Ugric database:  
analyzed text corpora and dictionaries  
for less described Ob-Ugric dialects
- 1st July 2014 – 30th June 2017
- 1–5 team members for the Yugan Khanty team

<http://oudb.gwi.uni-muenchen.de/>

# 1.2 Project Background



## OUIDB

Index | Sitemap | LMU-Homepage | Login

ENGLISH | РУССКИЙ | DEUTSCH

The **OUIDB Yugan Khanty corpus (YK)** contains texts from an unpublished collection by Zsófia Schön, partly together with Lyudmila Nikolaevna Kayukova. The transcription depicts spoken language by means of a broad phonematization using IPA characters. Only words or, respectively, sentences from native speakers of Khanty are glossed and included in the lexicon.

Yugan Khanty (2) optional dialect all genres  glossed only  audio only all translation: Search

sort by original title  annotated only

Original_Title	Dialect	English_Title	ID	Genre_Form	Genre_Content	glossed	annot.	Audio	Translation
е j њо jvəliyənyugan khanty o piseyən (ENK)	(YK)		1533	prose (pro)	Tales (tal)	glossed	-	Audio	
е mp tju:n to nap i ki (TMK)	yugan khanty (YK)	Old-dog-backtendoned-man (TMK)	1514	prose (pro)	Ethnographic Texts (eth)	glossed	-	Audio	en dt ru
е mp tju:n to nap i mi (AIK)	yugan khanty (YK)	Old-dog-backtendoned-woman (AIK)	1515	prose (pro)	Ethnographic Texts (eth)	glossed	annot.	Audio	en dt ru
е mp tju:n to nap i mi (SPK)	yugan khanty (YK)	Old-dog-backtendoned-woman (SPK)	1518	prose (pro)	Ethnographic Texts (eth)	glossed	-	Audio	en dt ru
i ttən əntə neripti (TMK)	yugan khanty (YK)	Why one shouldn't misbehave at night (TMK)	1507	mixed (mix)	Ethnographic Texts (eth)	glossed	-	Audio	en dt ru
i ttən əntə neripti (AJM)	yugan khanty (YK)	Why one shouldn't misbehave at night (AJM)	1523	prose (pro)	Tales (tal)	glossed	-	Audio	en dt ru
i ttən əntə neripti (AIK)	yugan khanty (YK)		1538	mixed (mix)	Tales (tal)	glossed	-	Audio	ru
i ttən əntə neripti (JFP)	yugan khanty (YK)	Why one shouldn't misbehave at night (JFP)	1542	mixed (mix)	Tales (tal)	glossed	-	Audio	en dt ru
oŋq bu tieŋkaliyən o piseyən 2. (TMK)	yugan khanty (YK)	The nuthatch and his older sister 2. (TMK)	1479	prose (pro)	Tales (tal)	glossed	-	Audio	en dt ru
oŋq bu tieŋkaliyən o piseyən 1. (TMK)	yugan khanty (YK)		1596	mixed (mix)	Tales (tal)	glossed	-	Audio	
o l m ə l	yugan khanty	The strange little	1526	mixed (mix)	Tales (tal)	glossed	-	Audio	en dt ru

= 30 texts

# 1.3 Project Background



Yugan Khanty • spoken in North-West Siberia, along the rivers Bolshoy Yugan, Malyy Yugan, Salym and Demyanka



# 1.4 Project Background



- Yugan Khanty
- ca. 900 native speakers
  - Eastern Dialect Group > Surgut Khanty
  - only “south” subdialect of Surgut Khanty
  - phonotactic: V1 + /y/ + V2 > V1:

(1) *pu:γəʃ ~ pu:ʃ*  
‘village’

(2) *nɬŋ kinʲtʲe:*  
*nɬŋ kinʲtʲe-ye*  
2sg than-2sg  
‘than you’

## 2. Fieldwork Recordings





## 2.1 Fieldwork Recordings



Collectors: • Ludmila Nikolaevna Kayukova  
• Zsófia Schön

Collection: • 29 texts, 1 h 16 min 15 sec

Genre: • tales  
• ethnographic texts  
• personal accounts

Form: • prose  
• dialogues  
• prose + dialogues

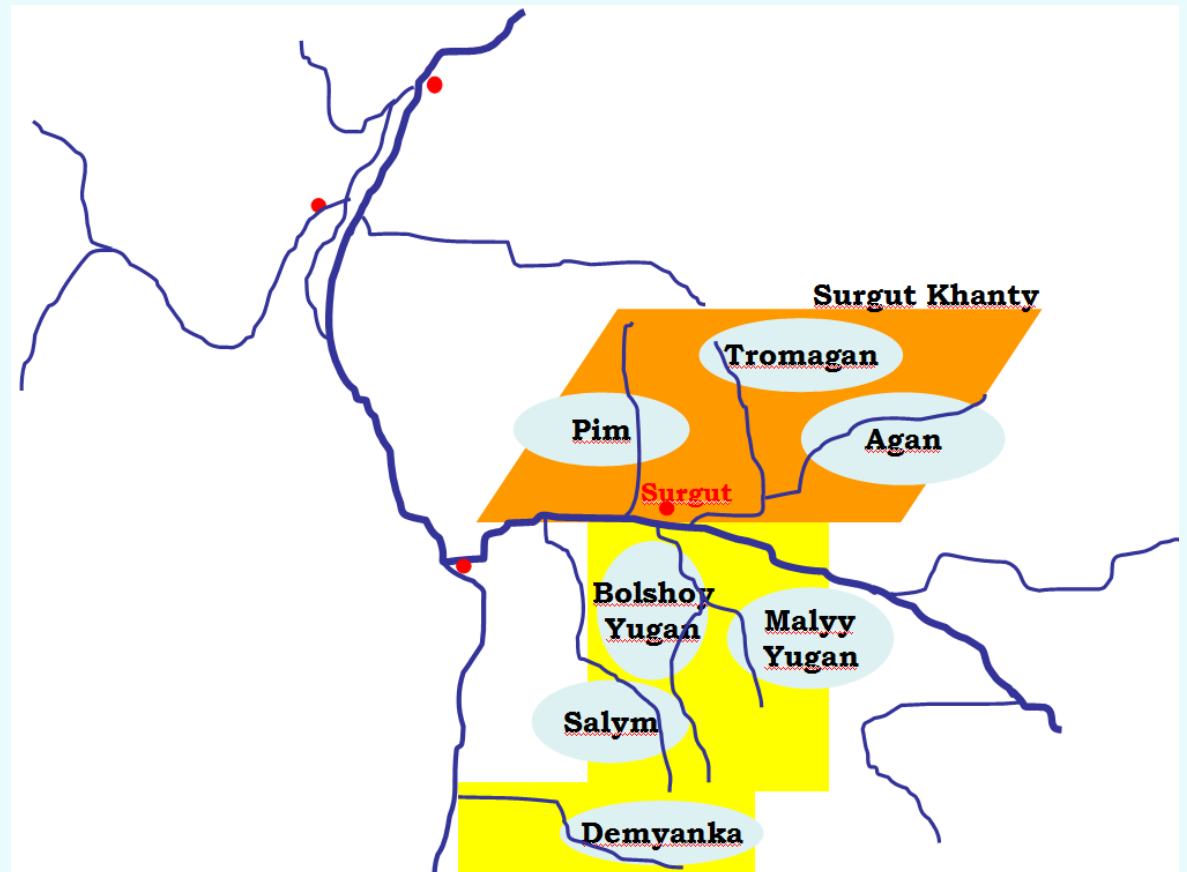


## 2.2 Fieldwork Recordings



Informants:

- 13 – 12 female and 1 male speaker
- Bolshoy Yugan
- Malyy Yugan
- Salym
- Demyanka





## 2.3 Fieldwork Recordings

- Equipment:
- 2010–2012 Olympus LS10
  - 2015–2017 Zoom H5
  - WAV 44.1kHz/16bit
  - stereo
  - built-in microphones
  - fixed position in front of the speaker



## 2.4 Fieldwork Recordings

Circumstances: • no laboratory circumstances



# 3. Corpus





## 3.1 Corpus

Contains: • 29 texts, 1 h 16 min 15 sec

Transcription: • IPA

- wide phonemic
- non-standardized, non-normalized
- spoken language
- speech disfluency
- elision, epenthesis

## 3.2 Corpus



- Workflow:
- segmentation and transcription in ELAN (29)
  - check, correction and on-the-fly Russian translation with L.N. Kayukova
  - glossing in FLEx (27)
  - translation into English (22), German (22) and Russian (23)
  - import of sound file, ELAN-Export and FLEx-Export in the database/user interface



# 3.3 Corpus

Data access: <http://oudb.gwi.uni-muenchen.de/>

- Metadata

Text + Metadata		Translation	Audio + Metadata	Glossed Text			
Original Title	Dialect	Informant	Genre Form	Genre Content	ID	glossed	Audio
i tten əntə nscripti (TMK)	yugan khanty (YK)	Kel'mina, Tat'yana Mikhaylovna (Kaymysova)	mixed (mix)	Ethnographic Texts (eth)	1507	by Schön, Zsófia	Audio
Text Source			Editor		Collector		
First publication Zsófia Schön (2016)			Schön, Zsófia; Kayukova, Lyudmila Nikolaevna		Schön, Zsófia (ZS)		
English Translation	German Translation	Russian Translation		Hungarian Translation			
"Why one shouldn't misbehave at night (TMK)" by Eckmann, Stefanie; Schön, Zsófia	"Wieso man nachts nicht unartig sein soll (TMK)" by Eckmann, Stefanie; Schön, Zsófia	"Почему ночью нельзя баловаться (TMK)" by Kayukova, Ludmila Nikolaevna; Snigirev, Yury; Schön, Zsófia		-			
Citation							
Schön, Zsófia 2016: OUIDB Yugan Khanty (2010–) Corpus. Text ID 1507. Ed. by Schön, Zsófia. <a href="http://www.oudb.gwi.uni-muenchen.de/?cit=1507">http://www.oudb.gwi.uni-muenchen.de/?cit=1507</a> (Accessed on 2017-09-27)							

### Audio Metadata:

Informant	Year of Birth (Informant)	Subdialect (Informant)	Collector	Date of Recording	Place of Recording	recorded by
Kel'mina, Tat'yana Mikhaylovna (Kaymysova)	1966	Yugan Khanty	Schön, Zsófia (ZS)	2012	Yurty Achimovy II	Schön, Zsófia





# 3.4 Corpus

Data access: <http://oudb.gwi.uni-muenchen.de/>

- Audio + Metadata

YK\_1507\_ZS\_ittenteneripti.wav (right click and 'save as' to download)

play all  
play next sentence

TMK: əj mət tɛ:tɲə ɔtə əj mətə pu:tɲə nʲitə i:mip i:kip βo ...  
 TMK: a kɛ:t ... no nʲitə i:mip i:kip  
 TMK: ɔtə i:mi:n tit kɛ:t sɛmʲa:n βistʲu:  
 foreign researcher: pi:pkit eto ...  
 TMK: i:mit i: ... i:mit i:kip  
 foreign researcher: i:kip  
 TMK: i:mip nʲit ... nʲitə

- Glossed Text

Export Mode  Audio Mode

1

#	TMK	:	#	əj mət tɛ:tɲə		ɔtə	əj	mətə	pu:tɲə	nʲitə	i:mip	i:kip	βo	#...#
				əj mət tɛ:tɲə		ɔtə	əj	mətə	pu:t-nə	nʲitə	i:mi-p	i:ki-p	βo	
				əj mətə	ɬɛ:tɲə+dial.var.eli.var.	ɔtə	əj	mətə	pu:t+dial.var.-nə	nʲitə+dial.var.	i:mi-p	i:ki-p	βo	βo
				once		adv	fil	cardnum	ipro	subs-infl:n	cardnum	subs-deriv:n>n	subs-deriv:n>n	v
				eh	r	one	some	village-LOC	four	wife-COLL	husband-COLL	live		

TMK: Once ehr in some village there liv... four wives and husbands.  
 TMK: Es ähm leb... einmal in irgendeinem Dorf vier Ehefrauen und Ehemänner.  
 TMK: Однажды ээ... в одной деревне четверо [человек] – женщины с мужьями – жили.



# 3.5 Corpus

- Data access:
- The Language Archive
  - 17 of the 29 Yugan Khanty texts
  - WAV, ELAN-, FLEx-Data combined
  - small metadata

The screenshot displays a software interface for audio analysis. At the top, there are two horizontal tracks showing audio waveforms. Below these is a detailed transcription table with multiple columns and rows, each representing a different layer of linguistic analysis. The table is organized as follows:

Layer	1	2	3	4	5	6	7
ref@ABC [1]	1515						
orth@ABC [27]	ej βəsyə məɾə ɔ:məs ij ne: ##						
word@ABC [190]	ej	βəsyə	mərə	ɔ:məs	ij	ne:	##
morph@ABC [197]	ej	βəsy	mar	ɔ:m	ij	ne:	
morph-var@ABC [66]		+fr. v		+[PS			
lemma@ABC [197]	ej	βəs	mar	ɔ:m	ej	ne:	
gloss@ABC [197]	once	well	time	sit	one	wom	
pos@ABC [197]	adv	ptcl	subs	v	card	subs	
ft-rus@ABC [27]	Однажды, ну, одна женщина долго сидела [и не спала].						
ft-hun@ABC [0]							
ft-eng@ABC [27]	Once a woman sat [awake] for a long time.						
ft-fin@ABC [0]							
ft-deu@ABC [27]	Einmal saß eine Frau lange Zeit [und war wach].						



## 3.6 Corpus

- Futur plans:
- orthographic transcription, literarization  
> book of Yugan Khanty tales
  - expansion of the corpus
  - search for new funding and projects



# Bibliography

- Чепреги, Марта. 2017. Сургутский диалект хантыйского языка. Ханты-Мансийск.
- Gugán, Katalin–Schön, Zsófia. 2018. 9.5 East Khanty. In: Bakró-Nagy, Marianne – Laakso, Johanna – Skribnik, Elena (eds) Oxford Guide to the Uralic Languages.
- Скрибник, Елена К.–Шён, Жофия–Янда, Гвен Ева–Визиорек, Аксель–Снигирёв, Юрий–Экманн, Штефани. 2017. Обско-Угорская база данных: текстовые корпуса и словари Обско-Угорских диалектов”. Ханты-Мансийск.
- Wisiolek, Axel–Schön, Zsófia. 2016. Obugric Database: Corpus and Lexicon Databases of Khanty and Mansi Dialects. In: Proceedings of the Second International Workshop on Computational Linguistics for Uralic Languages. Szeged.
- Wisiolek, Axel–Schön, Zsófia. 2017. Obugric Database: Corpus and Lexicon Databases of Khanty and Mansi Dialects. Acta Linguistica Academica 64,3. Budapest.

# *pe:sipe βoʃitəy*



<http://oudb.gwi.uni-muenchen.de/>

zsofia.schoen@gmail.com