

Predstavljanje Balkanološkog
instituta SANU na Seminaru za
jezičke tehnologije i resurse

Teodora Vuković

19.11.2015.

Tekući projekti i istraživanja

- Jezik, folklor i migracije na Balkanu

TIMSKA TERENSKA ISTRAŽIVANJA

BALKANOLOŠKOG INSTITUTA SANU
(GRAĐA ARHIVIRANA U DABIJU 1999–2015)

Timska terenska istraživanja Balkanološkog instituta SANU (građa arhivirana u DABIju 1999–2015)



Multimedija baza kulturnog nasleđa na Balkanu

- Internet prezentacija građe Balkanološkog instituta
- Moguća pretraživost pomoću interaktivne mape
- Sortirano po temama, koje su vezane za vremenski interval na snimku


Bunjevački rečnik


- Audio rečnik **realne upotrebe jezika**
- Polaziše su **terenski snimci** prikupljeni od strane istraživača sa Balkanološkog instituta tokom 2009. godine
- U sastav rečnika ulaze **lekseme koje se ne koriste u standardnom srpskom jeziku**, ili se koriste ne drugačiji način
- Definicije su preuzete iz **Rečnika bačkih Bunjevaca** Marka Peića i Grga Bačlije
- Obrada podataka pomoću programa za obradu zvuka (SounForge, Audacity)

Bunjevački rečnik

- **disnotor** *im m* 1. svinjokolj, vreme klanja svinja. 

Primeri:

1. Svako je za Božić u to vreme zaklo jedno manje svinjče, da Božić opremi. A pravi disnotor, barem naše familije, smo pravili kad je dečiji raspust. 

1. Kad uđeš unutra onda čestitaš disnotor, tako se čestita i onda, ovaj, ima, grliš nekog. 

Bunjevački korpus

- Model dijalekatskog korpusa bunjevačkog govora
- 50 000 tokena / 45 000 reči
- Normalizacija u 3 nivoa prema standardnom srpskom jeziku
 - Morfofonološke razlike
 - Sintaksičke razlike i prevodive razlike u leksici
 - Nprevodive razlike u leksici
- POS anotacija normalizovanog teksta
- Lematizacija normalizovanog teksta
- Pretraživanje u Corpus Workbenchu pomoću komandne linije

Bunjevački korpus

original	norm1	norm2	norm3
samo	samo	samo	samo
imam	imam	imam	imam
šotošku	šotošku	šotošku	šotošku*
,	,	,	,
pa	pa	pa	pa
kad	kad	kad	kad
ne	ne	ne	ne
sćam	sećam	sećam	sećam
se	se	se	se
ja	ja	ja	ja
na	na	***	***
to	to	toga	toga

ORIGINAL	POS1	LEM1	NORM1	POS2	LEM2	NORM2	NORM3
S	Sp	s	S	Sp	s	S	S
otim	Pd	otaj	tim	Pd	taj	tim	tim
siče	Nc	siča	seče	Vm	seći	seče	seče
snopove	Nc	snopova	snopove	Nc	snop	snopove	snopove
.	#	.	.	#	.	.	.
A	C	a	A	C	a	A	A
sitna	Af	sitan	sitna	Af	sitan	sitna	sitna
pliva	Nc	pliv	pleva	Nc	pleva	pleva	pleva
,	#	,	,	#	,	,	,
prvo	Rg	prvo	prvo	Rg	prvo	prvo	prvo
iđe	Vm	ići	ide	Vm	ići	ide	ide
slama	Vm	slamati	slama	Vm	slamati	slama	slama
iz	Sp	iz	iz	Sp	iz	iz	iz
doba	Nc	doba	doba	Nc	doba	dreša	dreša

Šokački rečnik

- Polazište je rečnik govora Bačkog Brega Marka Ivoševa
- Rečnik je potrebo urediti po leksikografskim standardima
- Konsultovan rečnik
- Postoji mogućnost za dodavanje multimedijalnih sadržaja
- Aplikacija se može primeniti i na druge rečnike

Šokački rečnik

Dodata	Autor	Leksema	Audio i Slika	Etimologija
09/18/2015	svetlana	aljina		

RB	Vrsta reči	Oblik	Fonetska transkripcija	Audio
1	Im ž	aljine		

RB	Definicija	Specijalna leksika	Slika	Audio	Video	Primer						
1	ženska jednodielna odeća; haljina				VIDEO	<table border="1"><thead><tr><th>RB</th><th>Primeri</th><th>Audio</th></tr></thead><tbody><tr><td></td><td></td><td></td></tr></tbody></table>	RB	Primeri	Audio			
RB	Primeri	Audio										
2	odeća uopšte				VIDEO	<table border="1"><thead><tr><th>RB</th><th>Primeri</th><th>Audio</th></tr></thead><tbody><tr><td>1</td><td>Jeste 1, mamu, sprimili te moje aljine, vrijeme je da se dižem?</td><td></td></tr></tbody></table>	RB	Primeri	Audio	1	Jeste 1, mamu, sprimili te moje aljine, vrijeme je da se dižem?	
RB	Primeri	Audio										
1	Jeste 1, mamu, sprimili te moje aljine, vrijeme je da se dižem?											

Šokački rečnik

Unesite vrednosti u polja

Leksema

Browse...

No file selected.

Browse...

No file selected.

Oblici reči

RB	Vrsta reči	Oblik	Fonetska transkripcija	Audio
----	------------	-------	------------------------	-------

Etimologija

Frazeologizmi

RB

Definicija

Primer

Sinonimi

RB

Sinonim

1

Antonimi

RB

Antonim

Dodaj red

Unesi reč



Konstruisanje narativa

- Projekat u saradnji sa Humbolt univerzitetom u Berlinu
- Upotreba terenskih snimaka iz **srpsko-mađarskih naselja** u okolini Pančeva
- **Transkripcija** u programu **EXMARaLDA**
- Cilj je **korpus narativa** sa anotiranim elementima narativa po značenju i funkciji
- Svrha korpusa **je izučavanje teksta kao niza konstrukcija** (po analogiji sa rečeničnom sintaksom), gde su sastavni elementi konstrukcija, zasebne tematske i funkcionalne celine u okviru teksta

Konstruisanje narativa

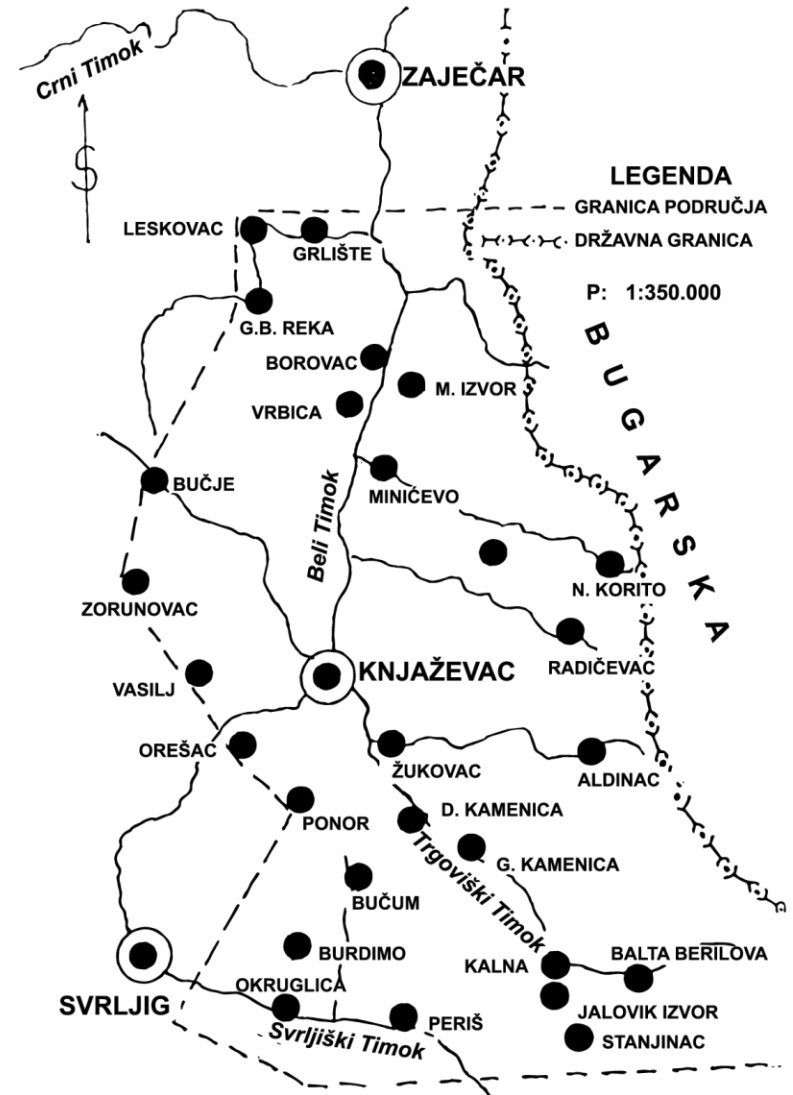
нагло умро. Тако да сам ја остала сама. Оставила ме и брат'а (плаче) ... Е тако. Али мој муж је Мађар био. Мој је муж Мађар био. Али он се прикључио мени и мојој вери. Тако да он је са

The screenshot displays an audio analysis interface. At the top, a waveform is shown with a time axis from 00:01 to 00:03. A vertical red line marks a point at approximately 00:02.6. Below the waveform is a control bar with a green arrow icon and the text "Append interval". To the right of this bar are two buttons labeled with asterisks and brackets. Below the control bar is a table with four columns and five rows. The first row shows time intervals: 0:00, 1 [00:00.2], 2 [00:00.2*], 3 [00:02.6], and 4. The second row contains the text "Али мој муж је мађар био. Али он се прикључио мени и мојој вери." The third row contains the labels "evaluacija" and "apstrakt". The fourth row repeats the text from the second row. The fifth row contains the labels "evaluacija" and "apstrakt".

0:00	1 [00:00.2]	2 [00:00.2*]	3 [00:02.6]	4
	Али мој муж је мађар био. Али он се прикључио мени и мојој вери.			
	evaluacija	apstrakt		
	Али мој муж је мађар био. Али он се прикључио мени и мојој вери.			
	evaluacija			
	apstrakt			

Timočki govori

- Terensko istraživanje u oktobru 2015
- Opštine Knjaževac, Zaječar i Svrljig



Timočki govori

- 150+ sati video i audio materijala
- 2000+ strana književnosti na dijalektu
- 100+ sagovornika
- 50+ sela

- Procena veličine – oko 1,5 milion reči

Timočki govori

- **Internet arhiv sa interaktivnom mapom,** putem koje se mogu pregledati materijali iz svakog od sela
- Video, audio, fotografije
- Obrada podataka pomoću programa za editovanje audio i video materijala - Audacity, Power Director

Timočki govori

- Transkripcija materijala u programu EXMARaLDA
- Izrada korpusa
- Morfološka i tematska anotacija
- Pretraživost po mestu, polu, obrazovanju, starosti...
- Izrada mape mikrodijalekata
- Kontrastivna analiza jezika u kontaktu

Hvala na pažnji!