

# “S kje pa si?”

## Metapodatki o regionalni pripadnosti uporabnikov družbenega omrežja Twitter

Konferenca “Slovenščina na spletu in v novih medijih”

Jaka Čibej<sup>1</sup> in Nikola Ljubešič<sup>2,3</sup>

<sup>1</sup>Filozofska fakulteta, Univerza v Ljubljani

<sup>2</sup>Filozofska fakulteta, Univerza v Zagrebu

<sup>3</sup>Institut “Jožef Stefan”

Ljubljana, 26. november 2015

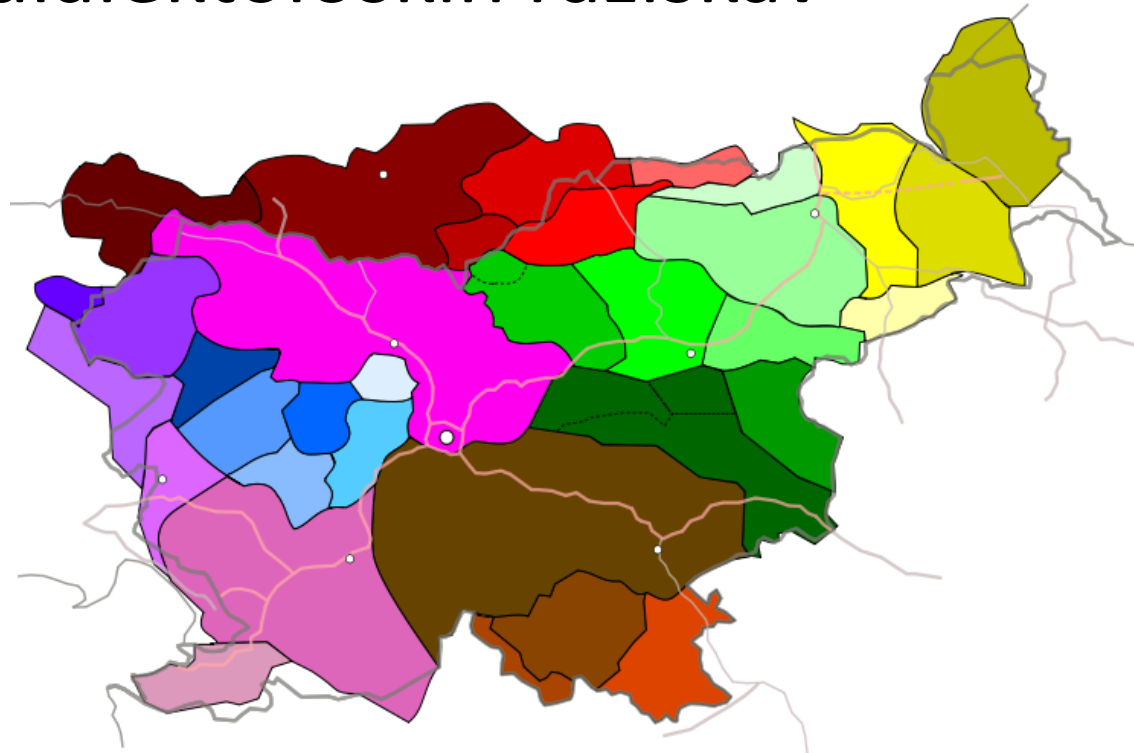
## Pregled predstavitve

1. Ozadje
2. Motivacija in sorodne raziskave
3. Dodajanje regionalnih metapodatkov v podkorpus tvitov JANES
4. Regionalni podkorpusi tvitov
5. Zaključek in prihodnje delo

# Ozadje

## Ozadje [1/4]

- 7 narečnih skupin, >30 narečij (Ramovš 1931)
- Veliko dialektoloških raziskav



## Ozadje [2/4]

- Narečja izumirajo (Kolarič 1954, Ramovš 1951)
- “Čista narečja” (idealnih govorcev)
  - standard *NORM (Non-Mobile Older Rural Male)* (Chamber & Trudgill 1994)
  - starost, izobrazba, spol, govornica staršev (Logar 1959)

## Ozadje [3/4]

- Razvoj narečij?
  - omejevanje na ruralno okolje (Sgall et al. 1992)
  - dialekte zamenjajo sociolekti (Sgall et al. 1992)
  - strnjevanje narečij (Niebaum & Macha 1999)
  - nova dialektizacija (Labov 1994)
- Ohranjanje narečij?
  - zavestno kultiviranje – Bavarska (Reichan 1999)
  - naklonjenost rabi narečja – Norveška (Jahr 1997)

## Ozadje [4/4]

- Sodobne korpusne dialektološke raziskave
  - nizozemščina → DynaSAND (Kunst & Wesseling 2010)
  - nordijski jeziki → Nordic Dialect Corpus (Johanessen et al. 2009)
  - angleščina → Freiburg Corpus of English Dialects (Hernández 2006)

# Motivacija in sorodne raziskave



## Motivacija

- Računalniško posredovana komunikacija
  - npr. Facebook, Twitter, Snapchat
  - vseprisotna (SURS)
    - družbena omrežja: 60 %
    - (video)telefoniranje: 41 %
- Spletni jezik / netspeak (Crystal 2001, Ueberwasser 2013, Erjavec & Fišer 2013)
  - mdr. regionalne jezikovne prvine

## Sorodne raziskave

- Regionalne jezikovne prvine v spletnem jeziku
  - do zdaj domena jezikovnih tehnologov
- Razvoj orodij za šumna besedila
  - avtomatska detekcija jezikov/različic jezika
    - arabščina (Cotterell & Callison-Burch 2014, Harrat et al. 2013)
    - ameriška angleščina (Eisenstein et al. 2015)
    - hrvaščina/srbščina/bosanščina/črnogorščina (Ljubešić & Kranjčić 2014)
  - oblikoskladenjsko označevanje
    - švicarska nemščina (Ruef & Ueberwasser 2013)
    - alzaščina (Bernhard & Ligozat 2013)
  - strojni prevajalniki
    - regionalna <-> sodobna standardna arabščina (Harrat et al. 2014)
    - dunajščina <-> standardna avstrijska nemščina (Haddow et al. 2013)

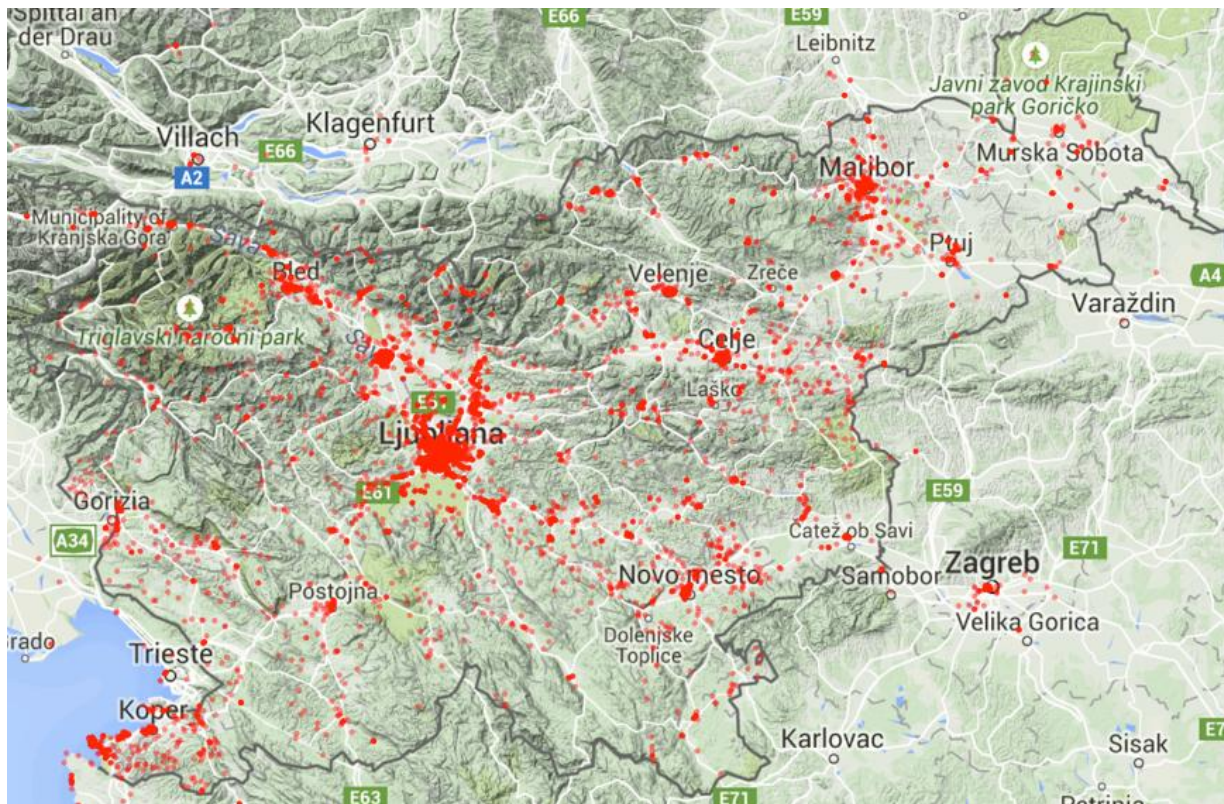
## Regionalne jezikovne prvine v spletni slovenščini

- Še neraziskane
- Korpus spletne slovenščine JANES
- Metapodatki o regionalni pripadnosti

# **Dodajanje regionalnih metapodatkov v podkorpus tvitov JANES**

## Zajem slovenskih tvitov z geolokacijo

- Orodje TweetCat (Ljubešič et al. 2014)
- začetek: januar 2015



## Izbor ustreznih avtorjev

- ~130.000 tвитov, ~1.600 avtorjev
  - izločitev računov organizacij na Twitterju
  - izločitev uporabnikov, ki niso vključeni v korpus JANES
- preostanek: 119.236 tвитov, 1.461 (zasebnih) avtorjev

## Razdelitev Slovenije na regije

- 9 regij → 7 narečnih skupin, Ljubljana, Maribor

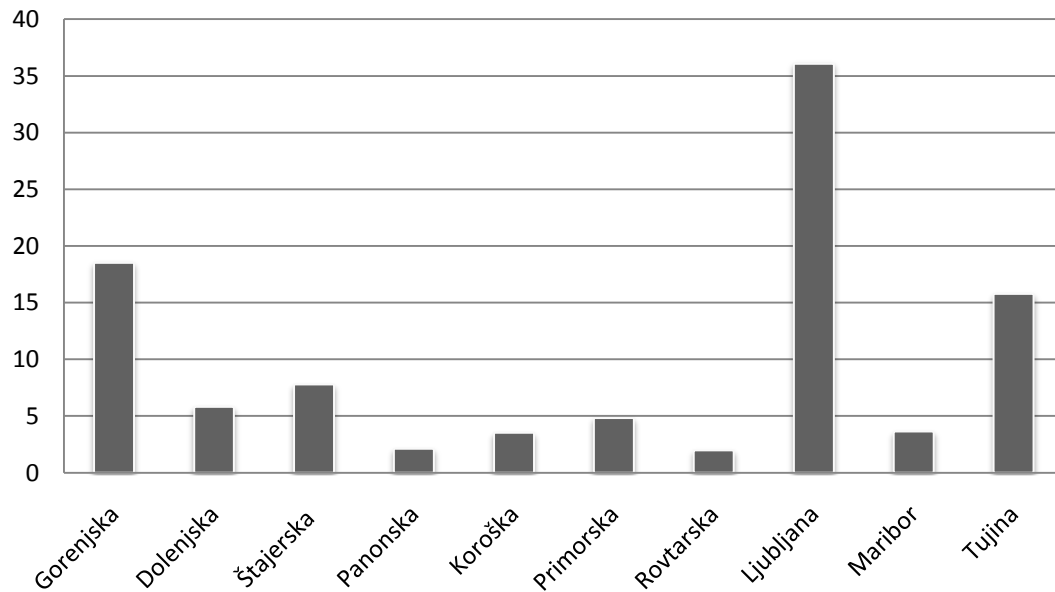


## Razvrščanje tvitov po regijah

- metoda metanja žarka (*ray-casting method*)



Deleži tvitov po regijah

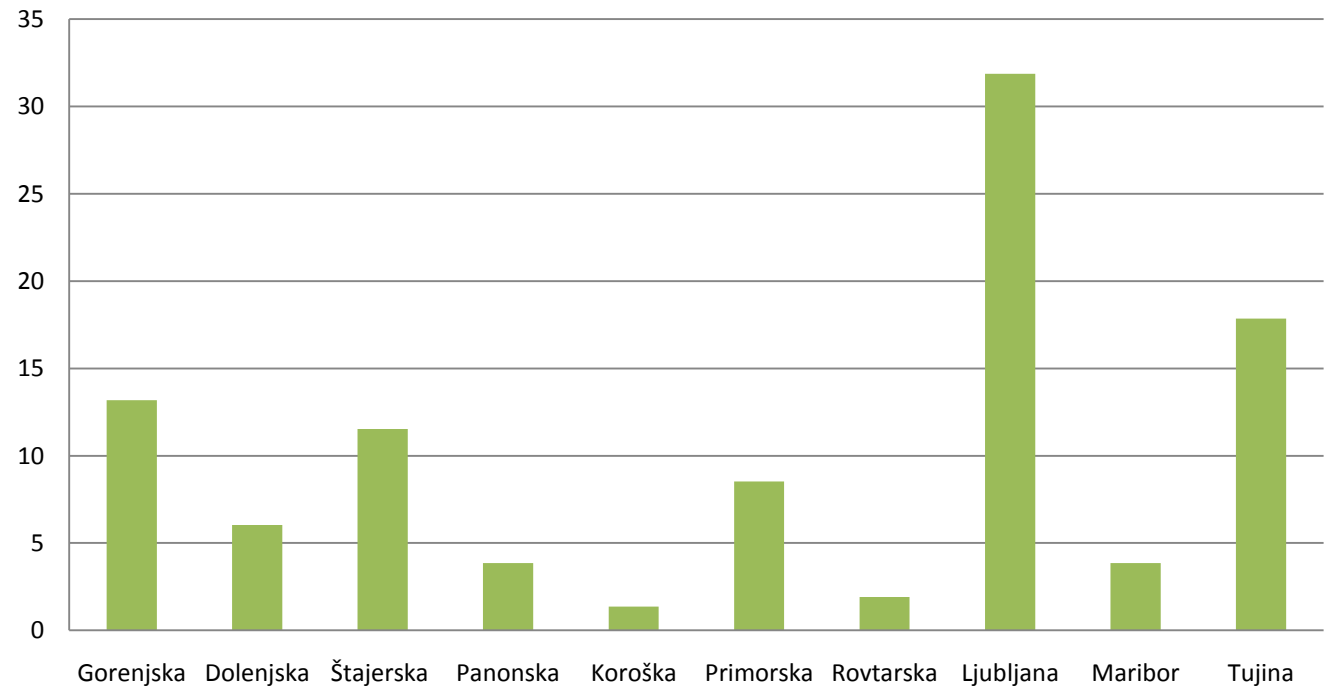




## Pripris metapodatkov

- Pogoji:
  - 90+ % tvitov iz ene regije
  - min. 3 tviti
- 364 avtorjev
  - povp. 74
  - med. 18
  - max. 1188

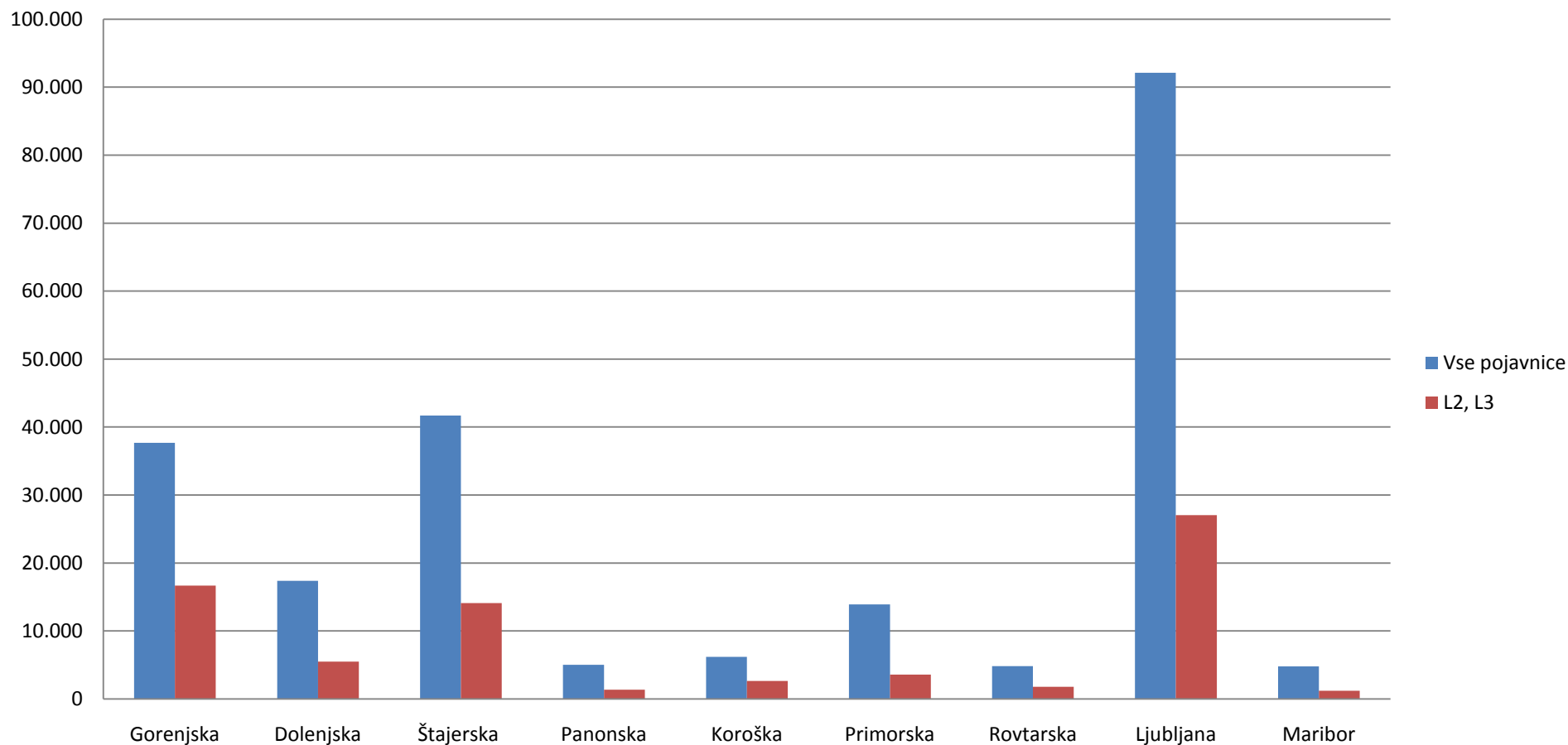
Deleži avtorjev po regijah



# Regionalni podkorpusi tvitov JANES

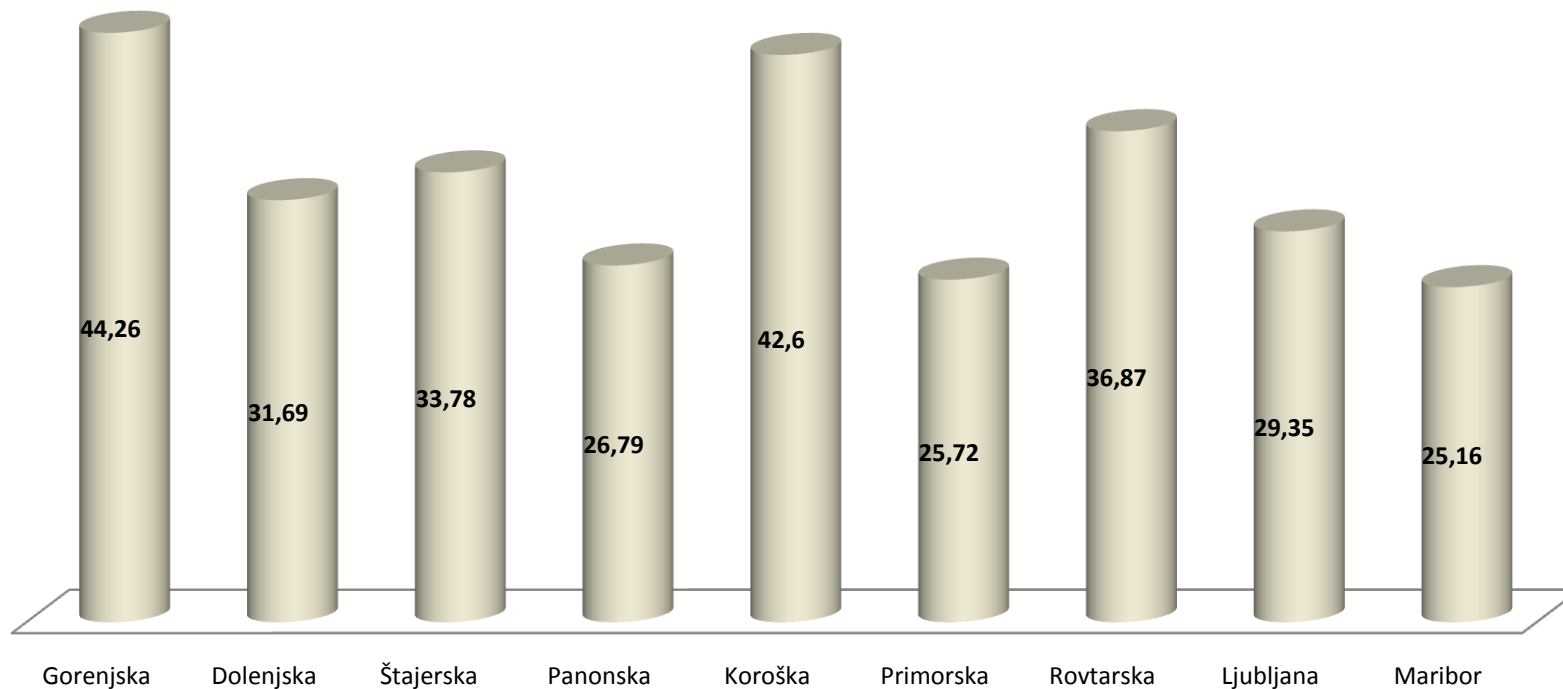
# Regionalni podkorpusi tvitov JANES

**Pojavnice v regionalnih podkorpusih**



# Regionalni podkorpusi tvitov JANES

Deleži L2/L3 v podkorporjih



# Zaključek in prihodnje delo

## Zaključek in prihodnje delo

- Postopek za avtomatski pripis metapodatkov o regionalni pripadnosti
- Prihodnje delo:
  - povečanje obsega podkorpusov (novi tviti)
  - izključitev premajnih podkorpusov
  - jezikoslovna analiza
  - primerjava s korpusom GOS
  - model za avtomatsko klasifikacijo regionalnih jezikovnih različic
    - primerjava z gručenjem

Hvala za pazornost.