

CLASSLA-web: Comparable Web Corpora of South Slavic Languages Enriched with Linguistic and Genre Annotation

Nikola Ljubešić and Taja Kuzman

Jožef Stefan Institute, University of Ljubljana, Institute of Contemporary History, Slovenia

CLASSLA-web Corpora

- First comparable corpus collection covering the entire language group
- The biggest general corpus for each South Slavic language
- Automatic genre and linguistic annotation

Web Corpora Sizes

⇒ 13 billion tokens

⇒ 26 million texts

Slovenian
2,153M tokens
4.06M texts

Bosnian
802M tokens
1.99M texts

Serbian
2,765M tokens
5.26M texts

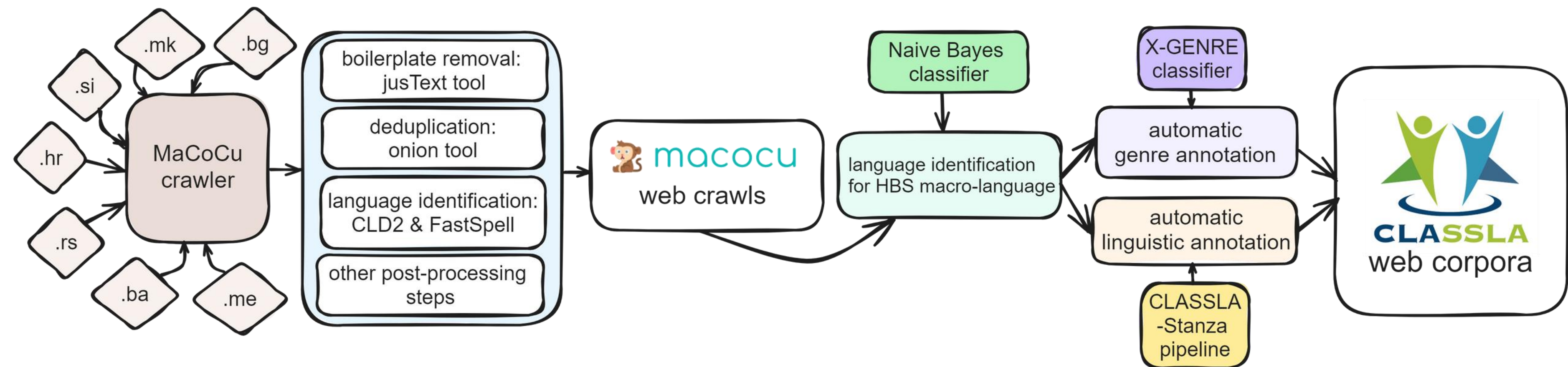
Bulgarian
3,918M tokens
7.46M texts

Croatian
2,576M tokens
5.42M texts

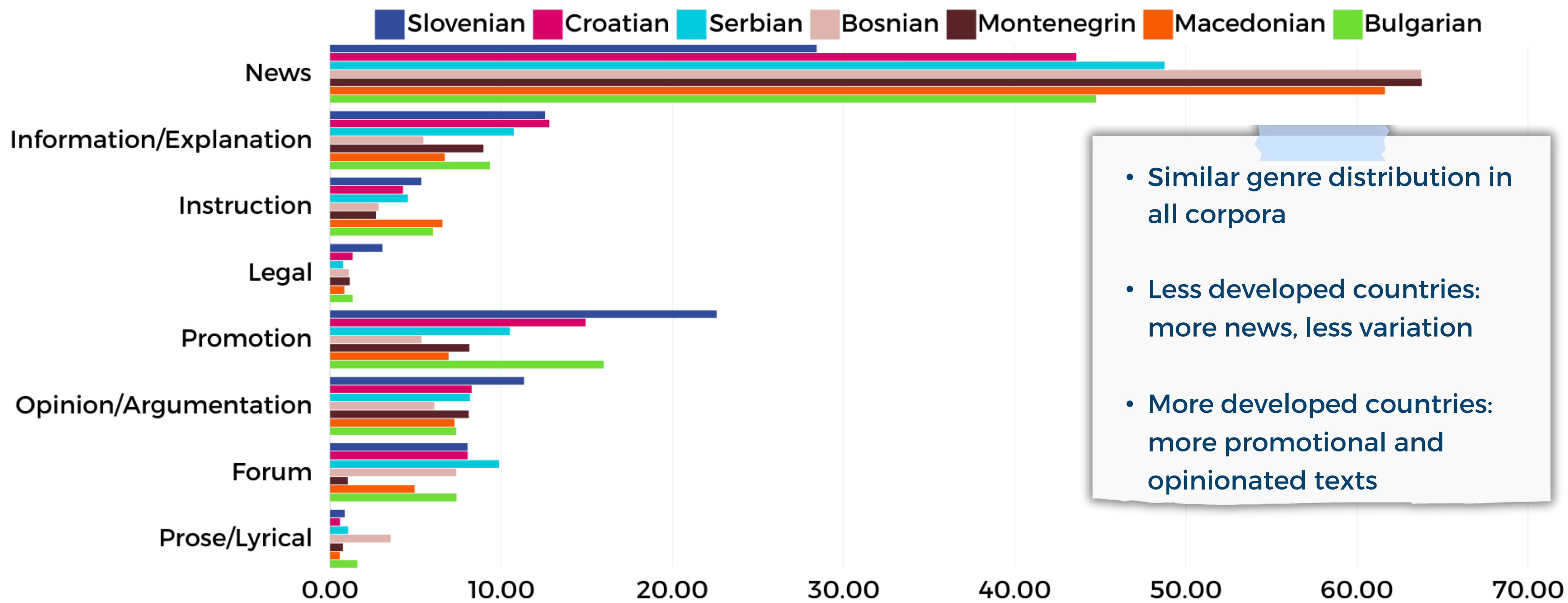
Montenegrin
177M tokens
0.40M texts

Macedonian
557M tokens
1.48M texts

Pipeline



Genres



- Similar genre distribution in all corpora
- Less developed countries: more news, less variation
- More developed countries: more promotional and opinionated texts

Plans for the Future

- Analyses of genre classifier's multilingual performance
- Series of workshops for wider linguistic community: CLASSLA-Express, visiting 5 South Slavic countries
- Iterative crawling to expand and update corpora: CLARIN.SI crawling infrastructure



Thank you