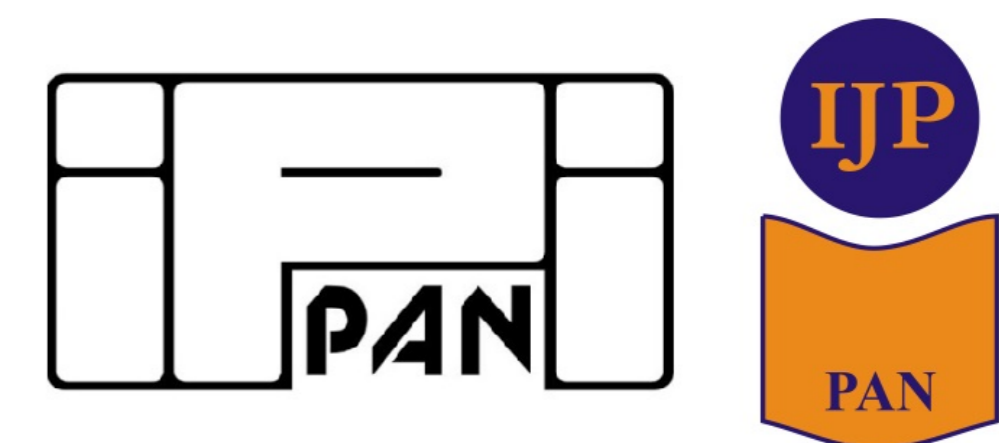# Morphosyntactic Annotation of Historical Texts
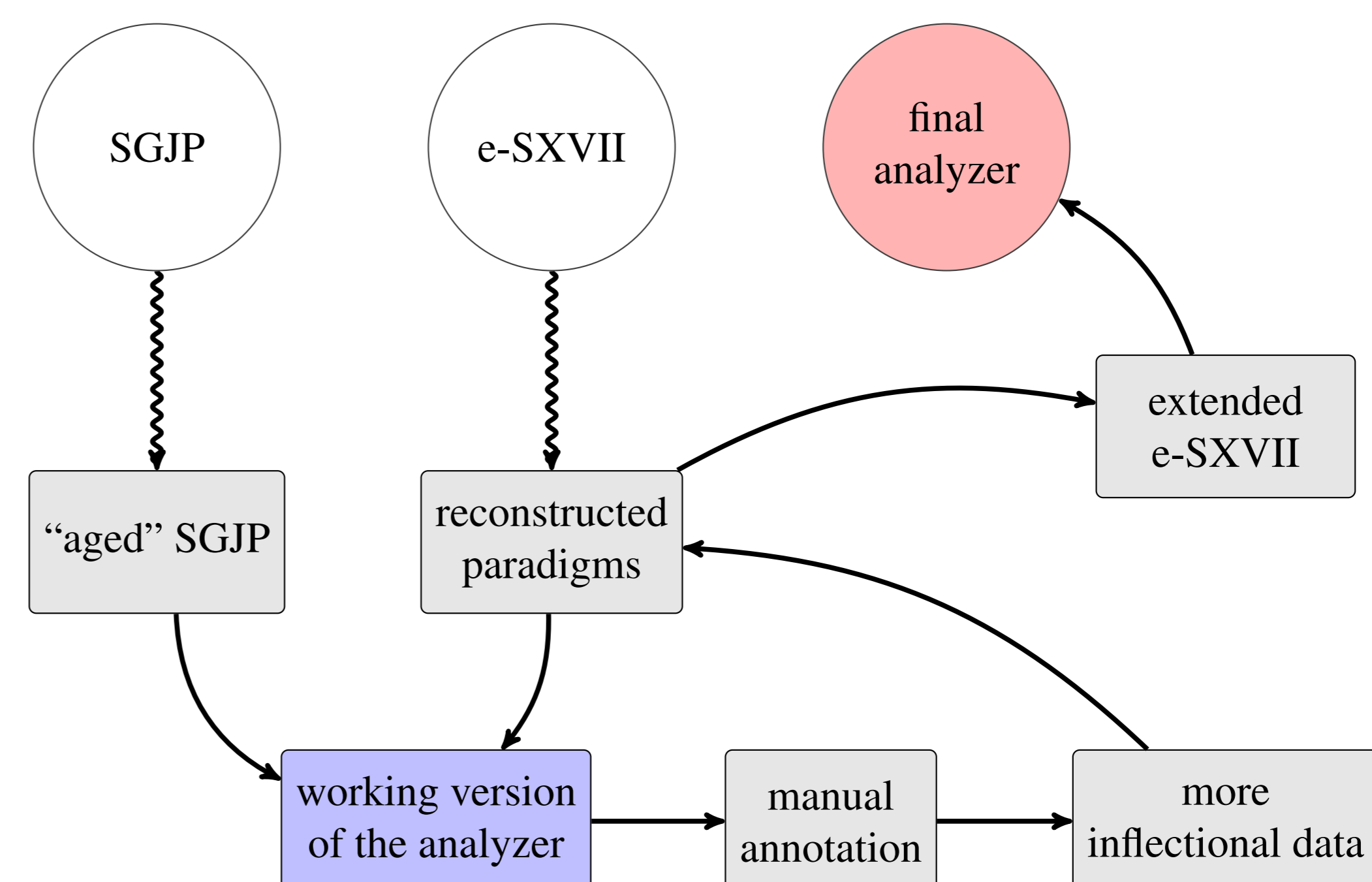
*The Making of the Baroque Corpus of Polish*

## W. Kieraś, D. Komosińska, E. Modrzejewski & M. Woliński

Institute of Computer Science & Institute of Polish Language
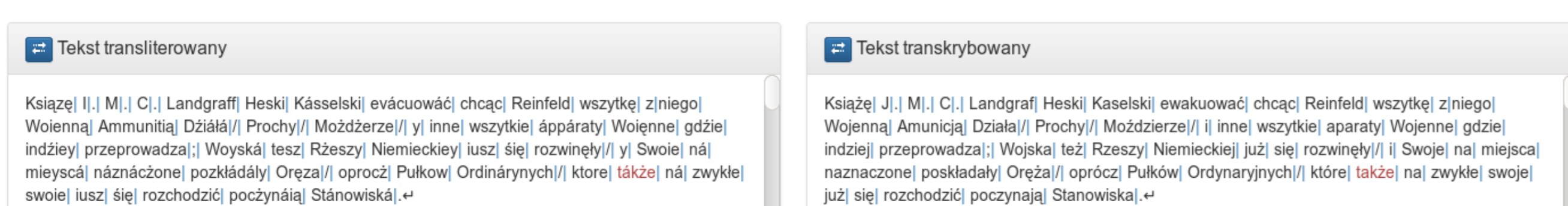Polish Academy of Sciences

**Abstract**

We present some technical issues concerning processing 17[th] & 18[th] century texts for the purpose of building a corpus of that period. We describe a chain of procedures leading from transliterated source texts to morphological annotation of text samples that was implemented for building the Baroque Corpus of Polish. The procedure consists of: automatic transliteration from original spelling to modern one, morphological analysis (including the construction of an inflectional dataset for Baroque Polish) and a tool for manual morphosyntactic annotation. The toolchain is being used to create a small manually validated subcorpus, which will serve as training data for a stochastic tagger. Then a larger corpus will be annotated automatically and made available via the Poliqarp corpus search tool.

## Automatic transcription

- The source texts of the Baroque Corpus of Polish divide into three types of editions: (1) original, (2) 19[th] century, and (3) contemporary.

- Types (1), (2) require automatic transcription to normalize the texts to make them better suited for morphological analysis.

- The converter contains nearly 4.000 substitution rules, based on regular expressions and the context of the sequence of characters' appearance.

- To check and correct the transcription, corpus annotators have a possibility to work on both transliterated and transcribed versions:



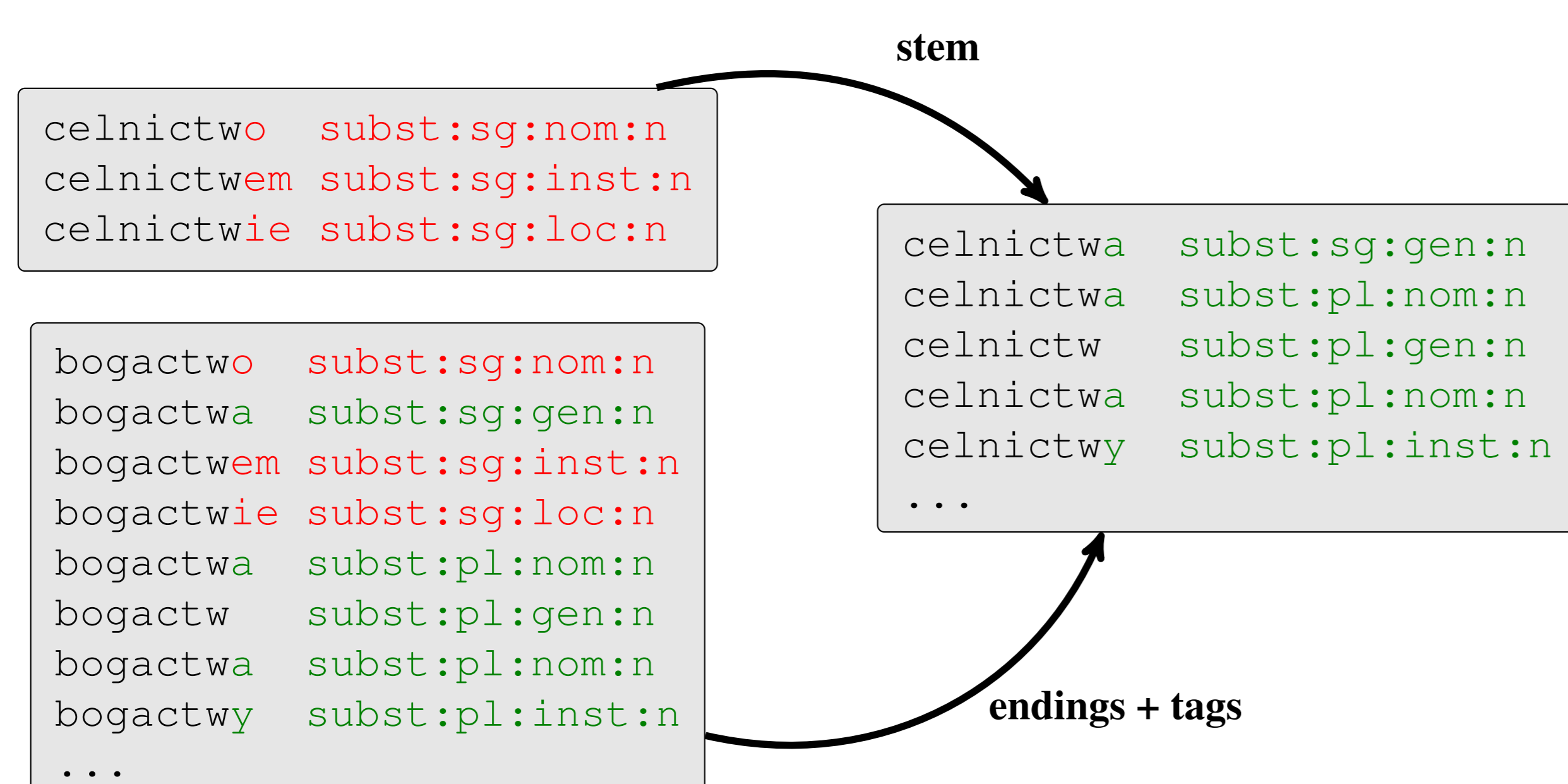A text sample manually transliterated (left) and then automatically transcribed (right)

- Automatic transcription normalizes the tokens, greatly reducing their types, and causes the number of tokens unknown to the automatic morphological analyzer (see below) to plummet to 5.4%:

|  | transliteration | transcription |
|---|---|---|
| token occurrences | 12,814,830 | 12,832,214 |
| token types | 646,410 | 476,733 |
| unrecognized token types | 75.55% | 48.08% |
| unrecognized token occurrences | 24.04% | 5.4% |

Number of tokens in the two representations of the corpus

## Morphological Analysis

- Automatic morphological analysis is performed using Morfeusz analyzer with a dictionary based on the following sources of data:

  - inflectional information from the Electronic Dictionary of 17[th] & 18[th] century Polish (e-SXVII, `http://sxvii.pl/`)

  - contemporary data of Grammatical Dictionary of Polish (SGJP) modified ("aged") to fit into the tagset for Baroque morphosyntax.

- Some automatically obtained extensions of inflectional paradigms from e-SXVII was also used. The extension procedure enriched the data set by over 210,000 new forms.
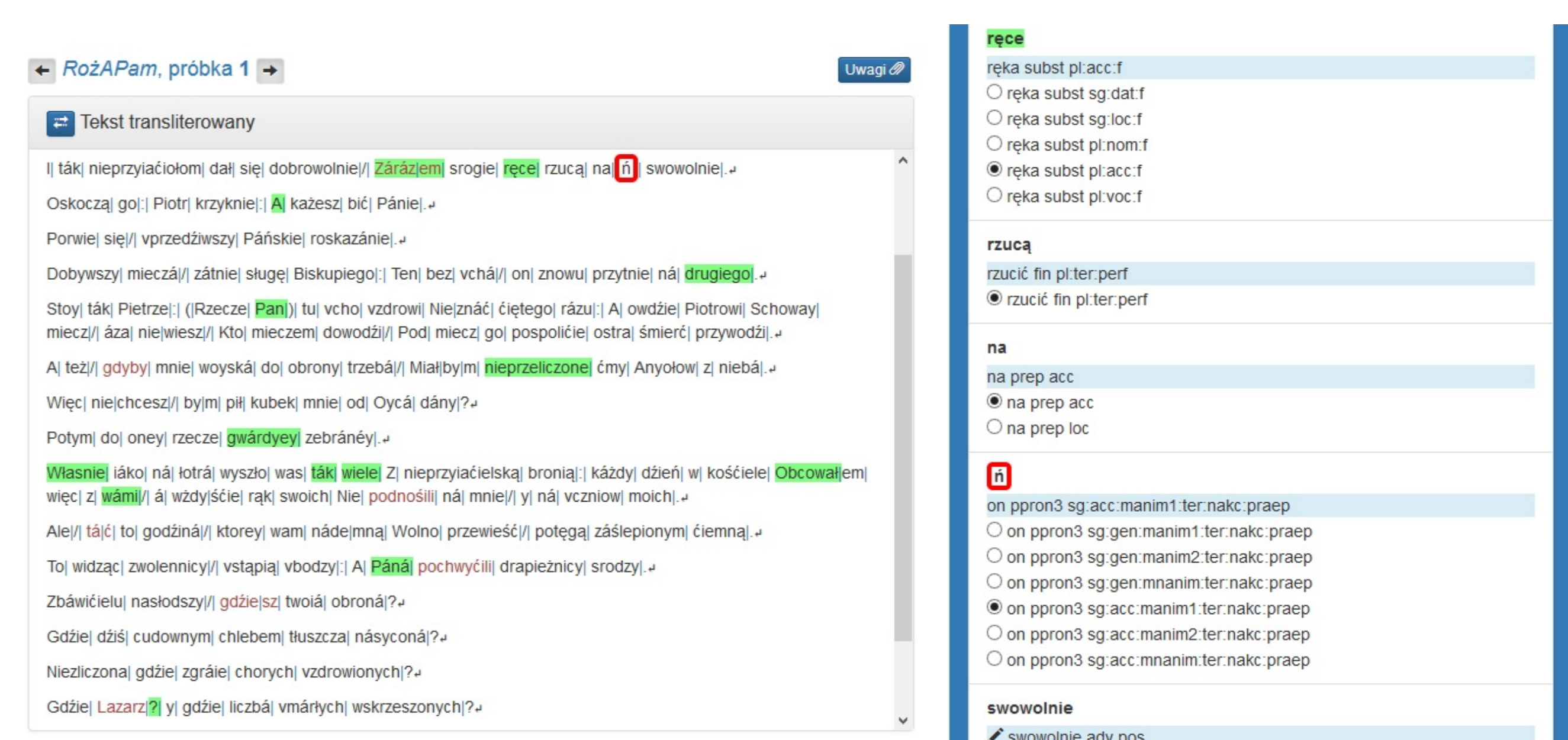


Example of partial reconstruction of inflectional paradigm from e-SXVII
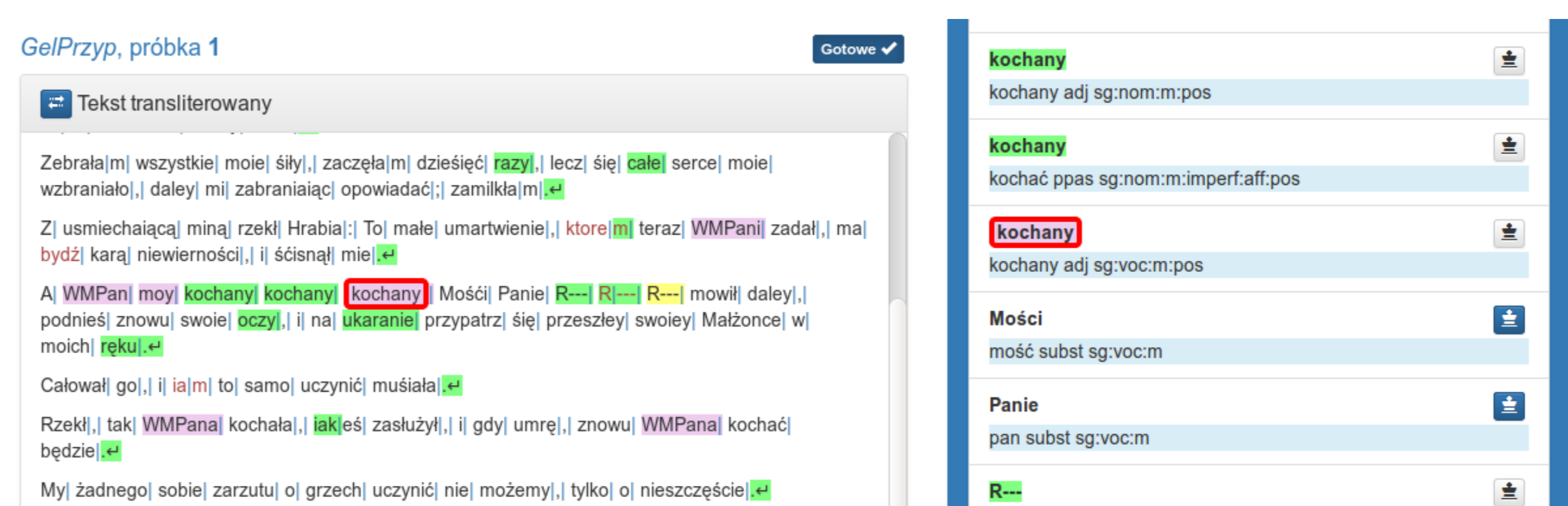
- "Aging" SGJP and reconstructing forms unnoted in e-SXVII allowed to reduce the number of tokens unrecognized by the analyzer in the corpus from 11% to 5,4%.

- Tokenization rules also needed to be adjusted. Baroque Polish, when compared to modern language, is characterized by wider use of joint spelling. Numerous clitics and particles are much more mobile and can be attached to a wider range of forms than nowadays.

- The analyzer is in the constant development based on the anotators' feedback. Its final version will be used for tagging the full corpus.



Workflow in the process of building the morphological analyzer

## Manual Annotation

- A corpus of 500,000 tokens of BCP will be manually annotated (currently ca. 150,000).

- All tokens marked as punctuation, foreign elements or structural markers were excluded for the purpose of the evaluation.

- The percentages of manual changes introduced by annotators are as follows:

  - modification of transcription – 2.66% of tokens, of tokenization – 1,59%,

  - morphological interpretation modified by the annotator – 8.6%.

- Each text sample is annotated independently by two annotators. In case of conflict the adjudicator has to resolve it by selecting one of the existing annotations or providing a new one.

- Currently annotators generate conflicts on 12.6% of tokens.



A sample being annotated as seen by an annotator in Anotatornia 2



Conflicts to be resolved as seen by an adjudicator

## Conclusions & further work

- The BCP is a work in progress since the manual annotation of the subcorpus is ongoing.

- All components involved in processing text samples are being constantly enhanced according to the feedback from annotators and adjudicators, but it seems that the tools are already quite stable.

- When the manually annotated corpus is ready, a stochastic tagger will be trained on the annotations and the rest of the BCP corpus will be automatically annotated. We also plan to check whether machine learning techniques can be used to train a better transcriber based on transcriptions manually corrected by the annotators.

- The presented toolchain will also be used to prepare a similarly organized corpus of 19[th] century Polish texts.

### Acknowledgements

5