# RUHR-UNIVERSITÄT BOCHUM

**RUB**

Fakultät für Philologie
Sprachwissenschaftliches Institut

# THE ANSELM PROJECT
## Tools for Automatic Analysis of a Parallel Corpus in Early New High German

**Marcel Bollmann, Stefanie Dipper, Julia Krasselt, Florian Petran**
**{bollmann,dipper,krasselt,petran}@linguistics.rub.de**

## The Project

In the Anselm project, we develop tools for the processing and automatic analysis of the medieval German treatise "St. Anselmi Fragen an Maria" (St. Anselm's questions to Mary), in collaboration with historical linguists from Ruhr-Universität Bochum.

Goals:
- Diplomatic transcription of all residual manuscripts and printings in German
- Annotation with bibliographic and linguistic information, e.g. POS tagging
- Alignment on paragraph, sentence, and word level
- Linguistic analysis of dialectal variation
- Publication of a digital edition

## The Corpus

In the texts, Anselm of Canterbury asks questions to the Virgin Mary concerning the Passion of Jesus Christ. She answers him in the form of longer monologues.

- 44 German manuscripts and prints
- Written between 14th to 16th centuries, in Early New High German (ENHG)
- Verse and prose versions of different lengths
- Comparable content, logical structure, and even (semi-)parallelity in sentence structure and wording
- No fixed spelling conventions
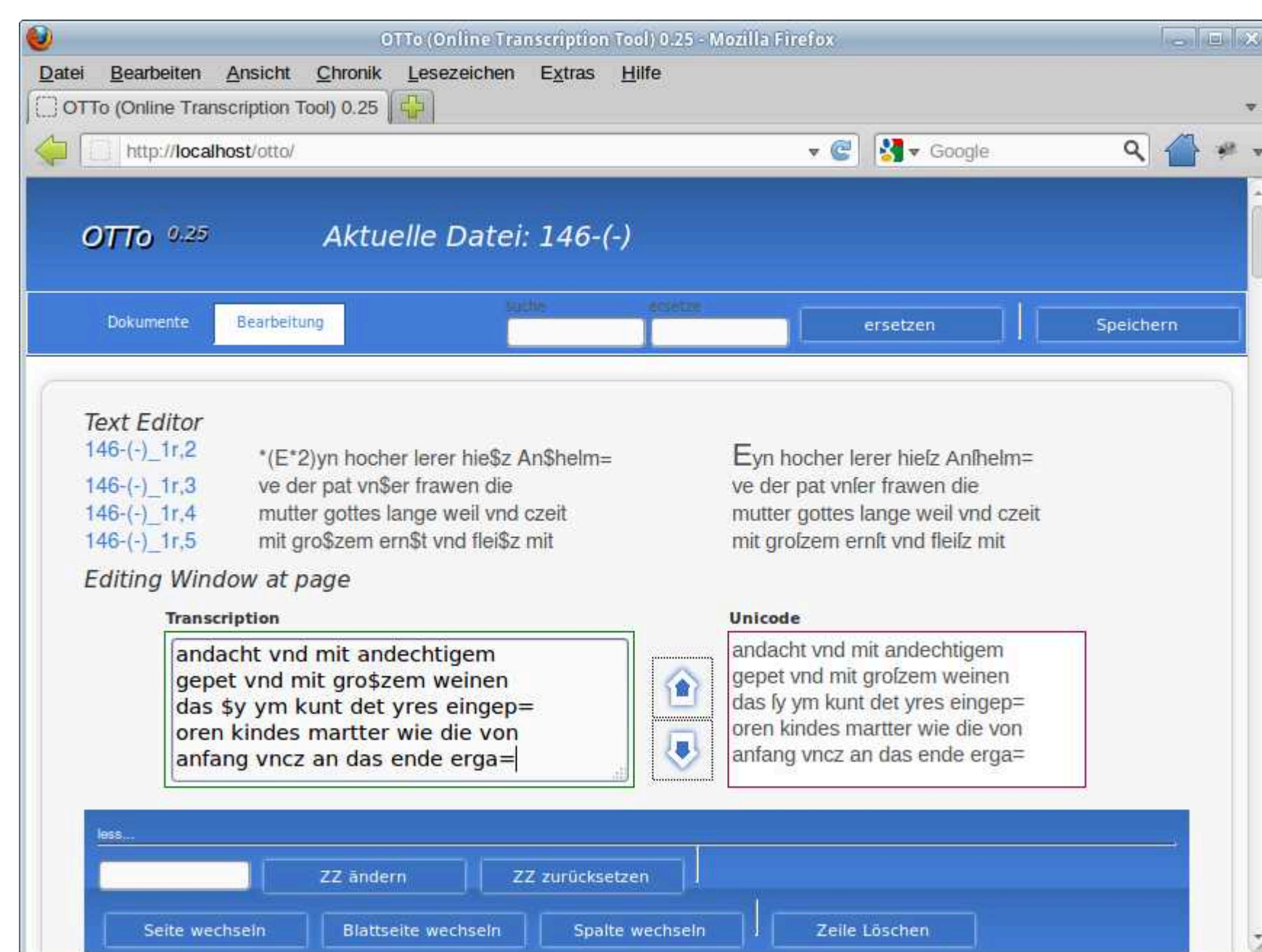- Dialectal variations in graphematics, phonology, morphology, and syntax

## The Tools

- Browser-based
  - No installation required by the user
  - Centrally maintained

**OTTo (Online Transcription Tool)**
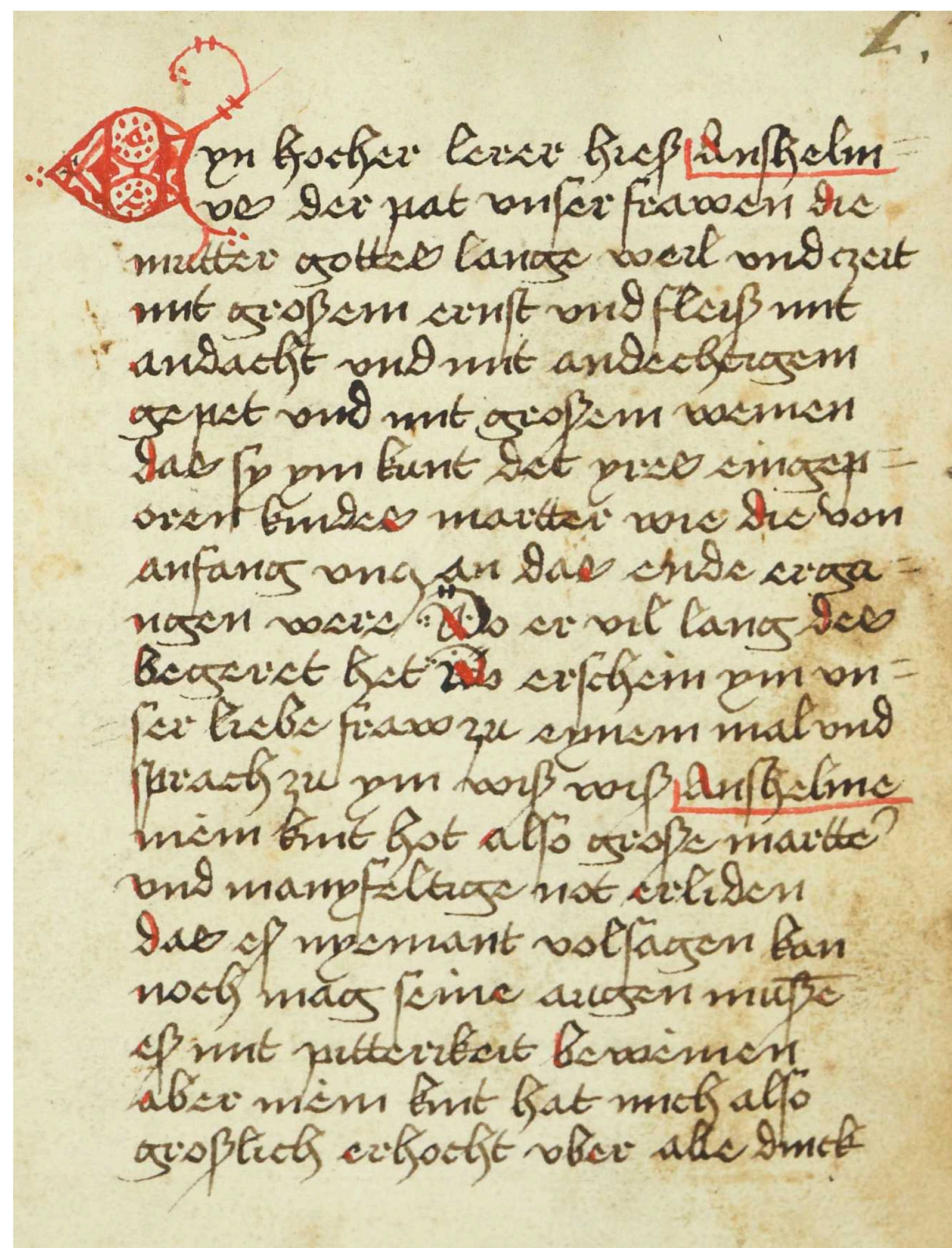
- Aids the diplomatic transcription of documents
- Special characters can be represented with user-definable escape sequences (e.g. '$' for 'ſ')
- Visual feedback: escape sequences are converted to Unicode characters and displayed instantly
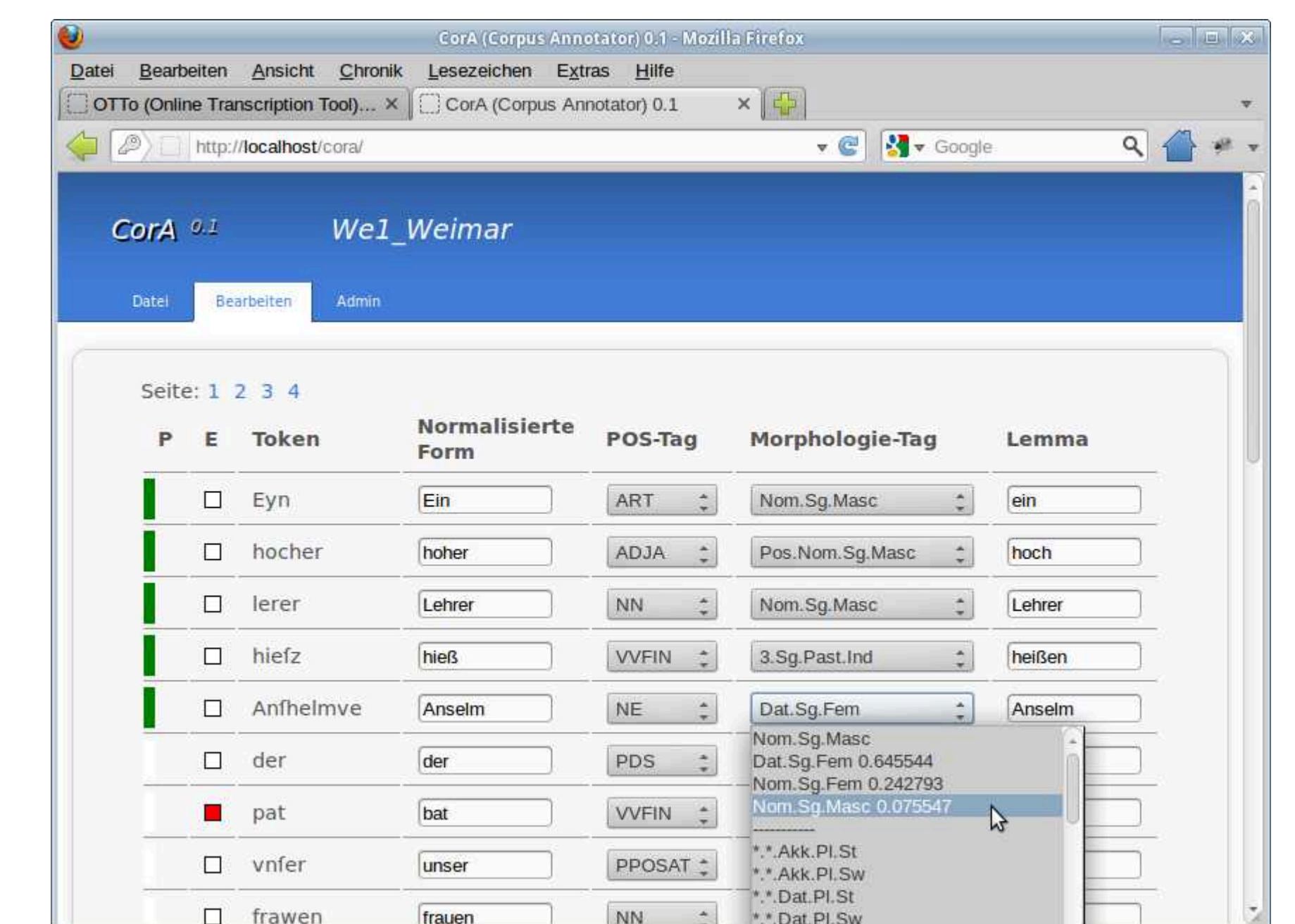
**CorA (Corpus Annotator)**

- Supports semi-automatic normalizations and annotations of various types


Screenshot of OTTo (Online Transcription Tool)


Sample from an Anselm manuscript (Upper German dialect), Weimar, Herzogin Anna Amalia Bibliothek, Cod. Oct. 4, fol.1r


Screenshot of CorA (Corpus Annotator)

## Transcription

- Diplomatic transcription with twofold collation
- Medieval texts use spelling conventions no longer in use today, e.g. long s ('ſ') or macrons ('ē')
- Graphematic peculiarities are potentially of interest for users of the corpus (e.g. for dialectal analysis) and should therefore be retained in transcriptions

## Alignment

As the texts are semi-parallel, alignment is often possible even at the word level.

Purpose:
- Projection of normalizations and annotations from one text over the whole collection
- Comparative linguistic and literary studies on the texts

Problems for traditional alignment methods:
- Missing or inserted sections
- Spelling differences hinder recognition of corresponding wordforms
- No sentence boundary markings

Solution:
- Extract a cognate dictionary
- Construct alignment sequences from closely occuring translation pairs
- Gradually expand the sequences to find the longest possible subsequence of alignment pairs
- Assign a confidence score and keep only the highest scoring sequence in case of conflict

Although it is a work in progress, this approach already performs better than traditional alignment methods based on heuristic text segmentation.

## Normalization

We map ENHG wordforms to New High German (NHG) wordforms, with the idea of ...

- applying already existing annotation tools for German (e.g. TreeTagger for POS tagging) to our data;
- facilitating the alignment process by standardizing the spelling; and
- measuring the distance between old and modern wordforms for purposes of linguistic analysis.

### Manual

- Manual normalization done by student assistants for evaluation purposes
- Two levels of normalization are distinguished:
  1) Modern wordform that is closest to the historical form ('vnſer frawen' → 'unser Frauen')
  2) Wordform adjusted with regard to inflection and semantics to form grammatical NHG phrases ('vnſer frawen' → 'unsere Frau')

### Automatic

- Idea: Learn characteristic spelling variations
- Character rewrite rules, sensitive to context

$$v \rightarrow u \,/\, \# \_ n$$

*(historic 'v' becomes modern 'u' between a word boundary and 'n', e.g. in the mapping 'vnd' → 'und')*

- Rules learned from an aligned Luther bible corpus
  - Based on the 1545 version and a modernized equivalent from 1892, both freely available
  - Alignment on paragraph, sentence, and word level
- Evaluation showed performance superior to a pure dictionary-based word substitution method
- Example mappings generated by our method:

| | | |
|---|---|---|
| etleich | → | etliche |
| vnse | → | unser |
| zitt | → | *zittern (instead of the correct 'zeit') |

- OTTo: Online Transcription Tool. http://www.linguistics.rub.de/otto/.
- Florian Petran (2012). *Aligning the un-alignable — a pilot study using a noisy corpus of nonstandardized, semi-parallel texts.* In: Alexander Gelbukh (ed.). Computational Linguistics and Intelligent Text Processing, Vol. 2. Springer: Berlin/Heidelberg.
- Marcel Bollmann, Florian Petran, and Stefanie Dipper (2011). *Applying Rule-Based Normalization to Different Types of Historical Texts — An Evaluation.* In: Proceedings of the 5th Language & Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics. Poznań, Poland.
- Marcel Bollmann, Florian Petran, and Stefanie Dipper (2011). *Rule-Based Normalization of Historical Texts.* In: Proceedings of the RANLP Workshop on Language Technologies for Digital Humanities and Cultural Heritage. Hissar, Bulgaria.

**DFG**

*An exalted teacher by the name of Anselm asked our lady, the mother of God, for a long while ...*

Eyn hocher lerer hieſz Anſhelmve der pat vnſer frawen die mutter gottes lange weil vnd czeit ...

Snct Anshelm batt vnser lieben frowen vom hymelrich lang zitt ...

*Saint Anselm asked our dear lady of the heavens for a long time ...*

Sample passage from two different manuscripts, illustrating word alignment