

## Motivation

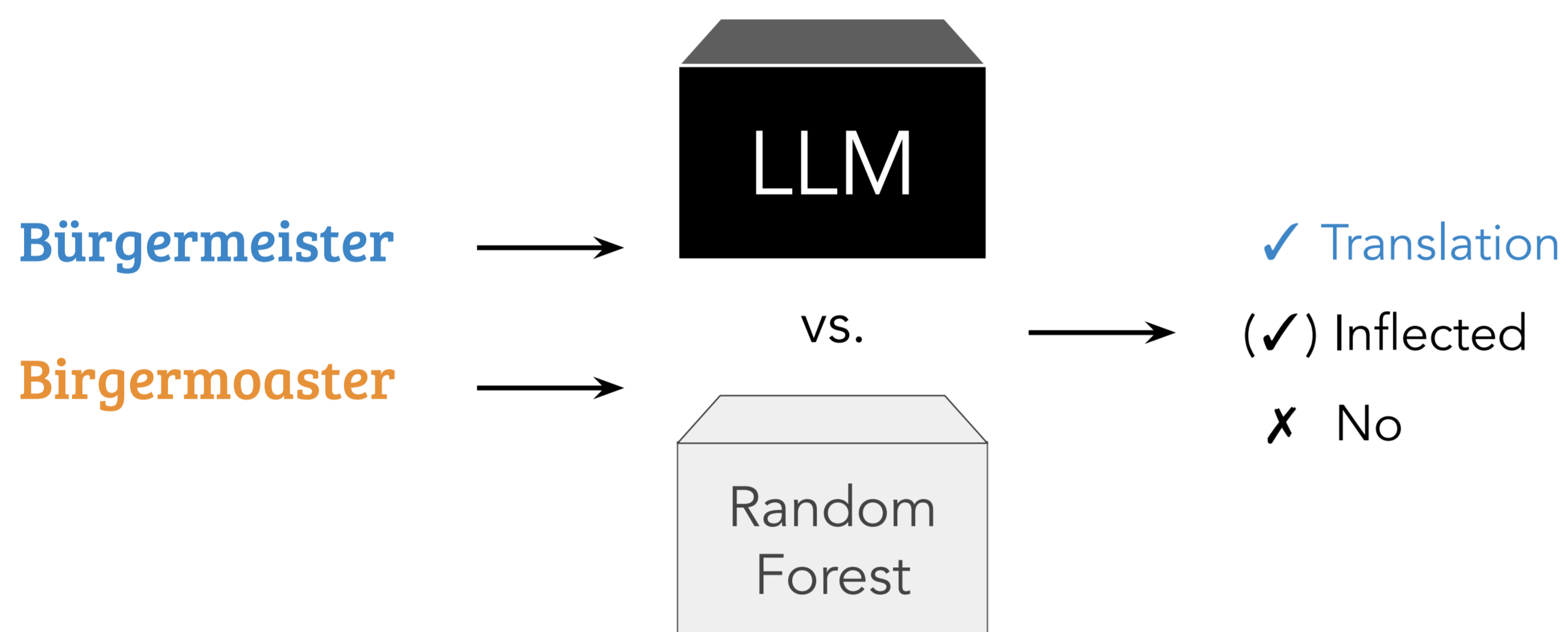
- Performance of NLP tools and LLMs depends on **how well language is represented** in training data.
- High-resource languages benefit, e.g., from **robust translation** and **retrieval** systems.
- Performance **degrades** under dialect spelling variation:

**Bürgermeister** (mayor) ↔ **Birgermoaster** **Bürgermaischer** **Birgermóaster** **Burgermoaster**  
**Byrgermoaster** **Bürgermoaster** **Birgermeister** **Birgermaster**

Bavarian spelling variations found in Wikipedia.

## Word-Pair Classification

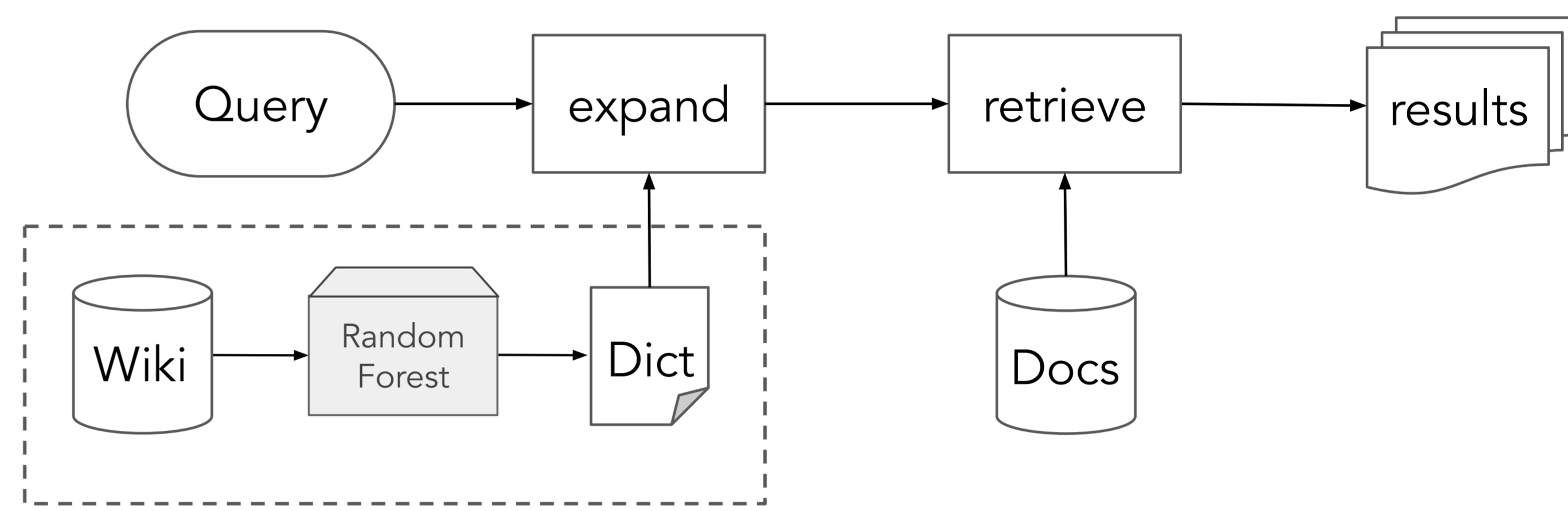
- Classify if pairs of German lemmas and dialect words are (direct or inflected) translations.
- LLMs **perform poorly** at identifying correct translations.<sup>1</sup>
- We evaluate a **statistical model** with **string similarity features**.



<sup>1</sup>Make Every Letter Count: Building Dialect Variation Dictionaries from Monolingual Corpora (Litschko et al., 2025)

## Cross-Dialect Information Retrieval

- **German-to-dialect retrieval**, spelling variations in documents.<sup>1</sup>
- Dialects: Riparian (ksh), Low German (nds), Alemannic (als), Rhine Franconian (pfl), Bavarian (bar).
- **Query expansion**: augment queries with spelling variations.



<sup>1</sup>Cross-Dialect Information Retrieval: Information Access in Low-Resource and High-Variance Languages (Litschko et al., 2025)

## String Similarity Features

Our string similarity features are adopted from prior work<sup>1</sup> and were originally developed for classifying words as cognates or false friends.

$$\text{DICE}(x, y) = \frac{2 \times \text{Intersection}}{\text{Union}}$$

rg, ge, er, rm, st  
Bü, ür, rg, ge, er, Bi, ir, rg, ge, er, rm, me, ei, is, st, te, rm, mo, oa, as, st

$$\text{LCSR}(x, y) = \text{B ü r g e r m e i s t e r} / \text{B i r g e r m o a s t e r}$$

### Set-based String Similarity

- **DICE** coefficient between ngrams(x) and ngrams(y):
  - Shared character bigrams, trigrams,<sup>2</sup> and “extended trigrams” (**XDICE**).<sup>3</sup>
- **XXDICE** incorporates positional information.<sup>3</sup>

### Sequence-based String Similarity

- Length of longest common prefix (**PREFIX**)
- Longest common subsequence ratio (**LCSR**;<sup>4</sup> uni-, bi-, trigram sequences)
- Length-normalized edit distance (**NED**;<sup>5,6</sup> uni-, bi-, trigram sequences)
- Edit distance between phonetic codes (**cologne phonetics**)<sup>7</sup>

Code and data:



<sup>1</sup> Automatic identification of cognates and false friends in French and English (Inkpen et al., 2005)

<sup>2</sup> The use of an association measure based on character structure to identify semantically related pairs of words and document titles (Adams and Boreham, 1974)

<sup>3</sup> Word-pair extraction for lexicography (Brew et al., 1996)

<sup>4</sup> Bitext maps and alignment via pattern recognition (Milamed, 1999)

<sup>5</sup> Binary codes capable of correcting deletions, insertions and reversals (Levenshtein, 1966)

<sup>6</sup> Identification of confusable drug names: A new approach and evaluation methodology (Kondrak and Dorr, 2004)

<sup>7</sup> Die Kölner Phonetik - Ein Verfahren zur Identifizierung von Personennamen auf der Grundlage der Gestaltanalyse (Postel, 1969)

## Results

