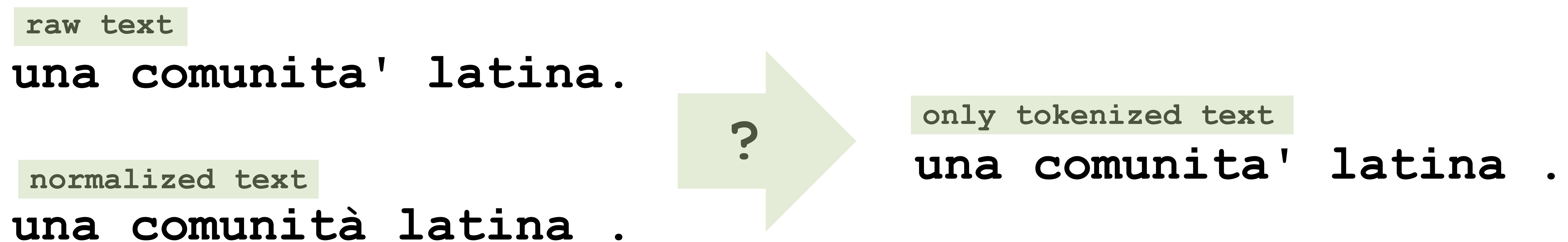


A Text Denormalization Algorithm

Producing Training Data for Text Segmentation

Problem



Solution

1. align Levenshtein matches



2. co-align heuristically



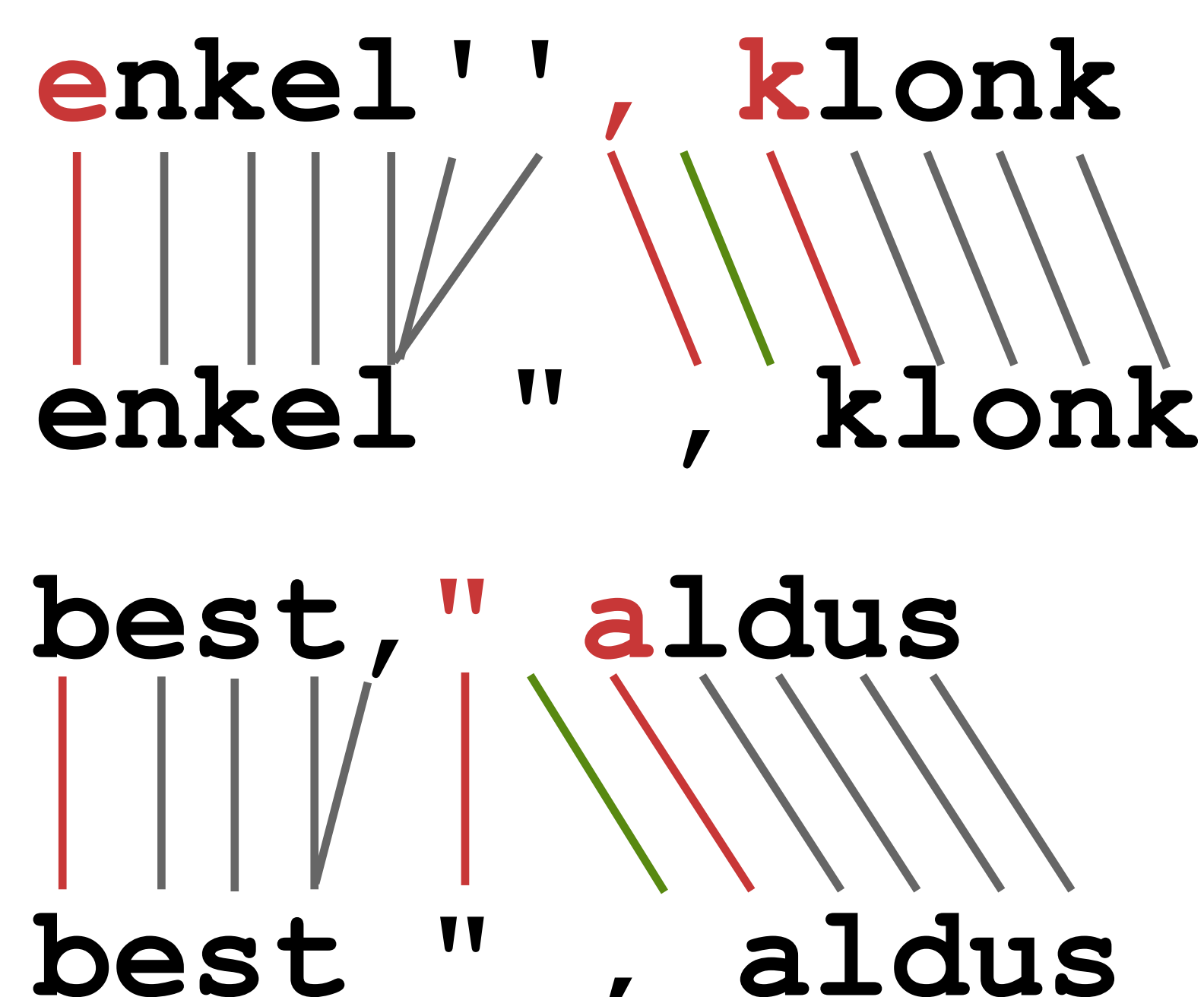
3. transfer boundaries



What works



What doesn't work yet



Conclusion

- Algorithm transfers token boundaries from normalized to raw text
- Levenshtein algorithm provides principled starting point
- Some normalizations require additional heuristics
- Application: train statistical text segmentation tools (Evang et al. 2013)