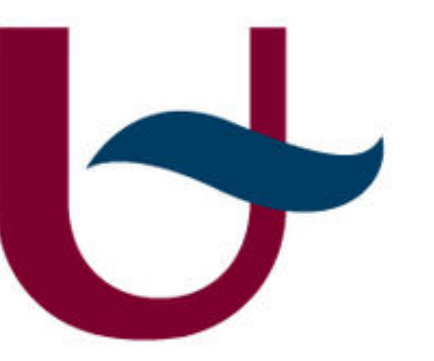


CLINICAL SPELLING CORRECTION FOR ENGLISH AND DUTCH WITH WORD AND CHARACTER N-GRAM EMBEDDINGS



Universiteit Antwerpen

{ PIETER FIVEZ, SIMON ŠUSTER, STÉPHAN TULKENS AND WALTER DAELEMANS } CLIPS, UNIVERSITY OF ANTWERP

Universiteit Antwerpen

CONTRIBUTION

- ✓ **CONTEXT-SENSITIVE** spelling correction for non-word misspellings in clinical text
- ✓ uses **WORD AND CHARACTER N-GRAM EMBEDDINGS** to estimate semantic fit
- ✓ **OUTPERFORMS** a state-of-the-art noisy channel model for **ENGLISH**
- ✓ **TRANSFERS** well to **DUTCH**

MOTIVATION

CLINICAL TEXT is characterized by

- high spelling **error** rate (mostly non-word)
- **variable** lexical characteristics
- **limited** accessible data

Applying the noisy channel model directly to this domain raises **THREE ISSUES**

1. Handling context-specificity

iron deficiency due to enemia → anemia
fluid injected with enemia → enema

2. Handling pseudo-random typos

patient not aware → are (correct: *aware*)

3. Handling imbalanced or highly variable training data

→ skewed corpus frequencies of domain-specific terms

ACKNOWLEDGEMENTS

This research was carried out in the framework of the Accumulate IWT SBO project, funded by the government agency for Innovation by Science and Technology (IWT). We also like to thank Elyne Scheurwegs for preparing the Dutch data.

REFERENCES

- [1] Bojanowski, P., E. Grave, A. Joulin, and T. Mikolov (2016). Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*.
- [2] Lai, K. H., M. Topaz, F. R. Goss, and L. Zhou (2015). Automated misspelling detection and correction in clinical free-text records. *Journal of Biomedical Informatics* 55, 188–195.

CONCEPT

Rank misspelling replacement candidates according to **SEMANTIC FIT**

→ extract salient clues from **misspelling context**

iron deficiency due to enemia → anemia

fluid injected with enemia → enema

→ **cosine similarity** between misspelling context and candidate in embedding space

EMBEDDING MODEL: FACEBOOK FASTTEXT [1]

- Word representations composed from word vectors and character n-gram vectors
- constructs vectors for **unobserved** word types
- more **generalized** representations
- incorporates distributional semantic information from **subword** level

DEVELOPMENT

CORPUS CREATION

- 12,740 instances from the MIMIC-III corpus
- transforming randomly sampled in-vocabulary words with random Damerau-Levenshtein edit operations
- 80% 1 edit operation: *anemia* → *enemia*
- 20% 2 edit operations: *anemia* → *enemea*

EXPERIMENTAL PARAMETERS

1. Vector composition functions

→ e.g. addition, multiplication

2. Context weighting

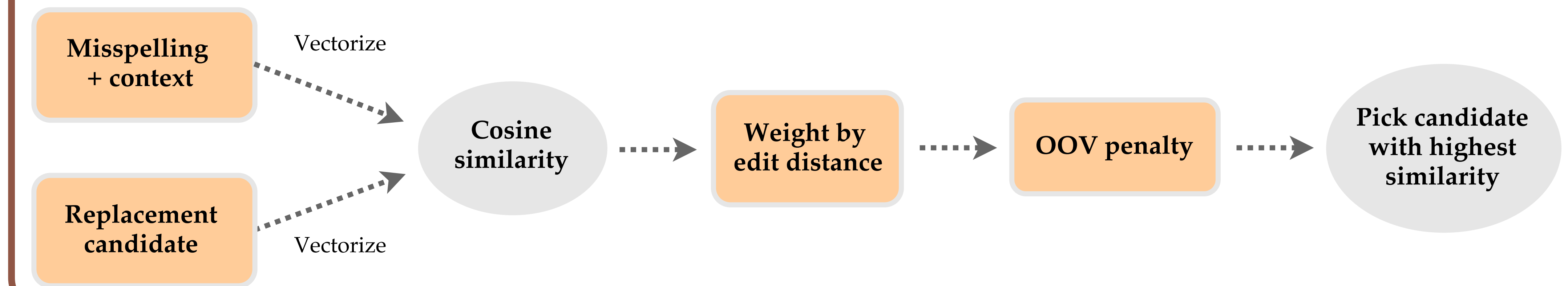
→ e.g. left vs. right window, window size

3. Orthographic and phonetic edit distance

4. Out-of-vocabulary (OOV) penalty

→ for candidate vectors constructed solely from character n-gram vectors

RANKING MODEL ARCHITECTURE



TEST SETTING

MODELS

CLINSPELL: our context-sensitive clinical spelling correction method

NOISY CHANNEL: our implementation of the state-of-the-art noisy channel method for clinical spelling correction by [2]

CORPORA

ENGLISH: 555 annotated observed non-word misspellings from MIMIC-III

DUTCH: 266 annotated observed non-word misspellings from heterogeneous clinical data from the University Hospital of Antwerp (UZA)

TEST RESULTS

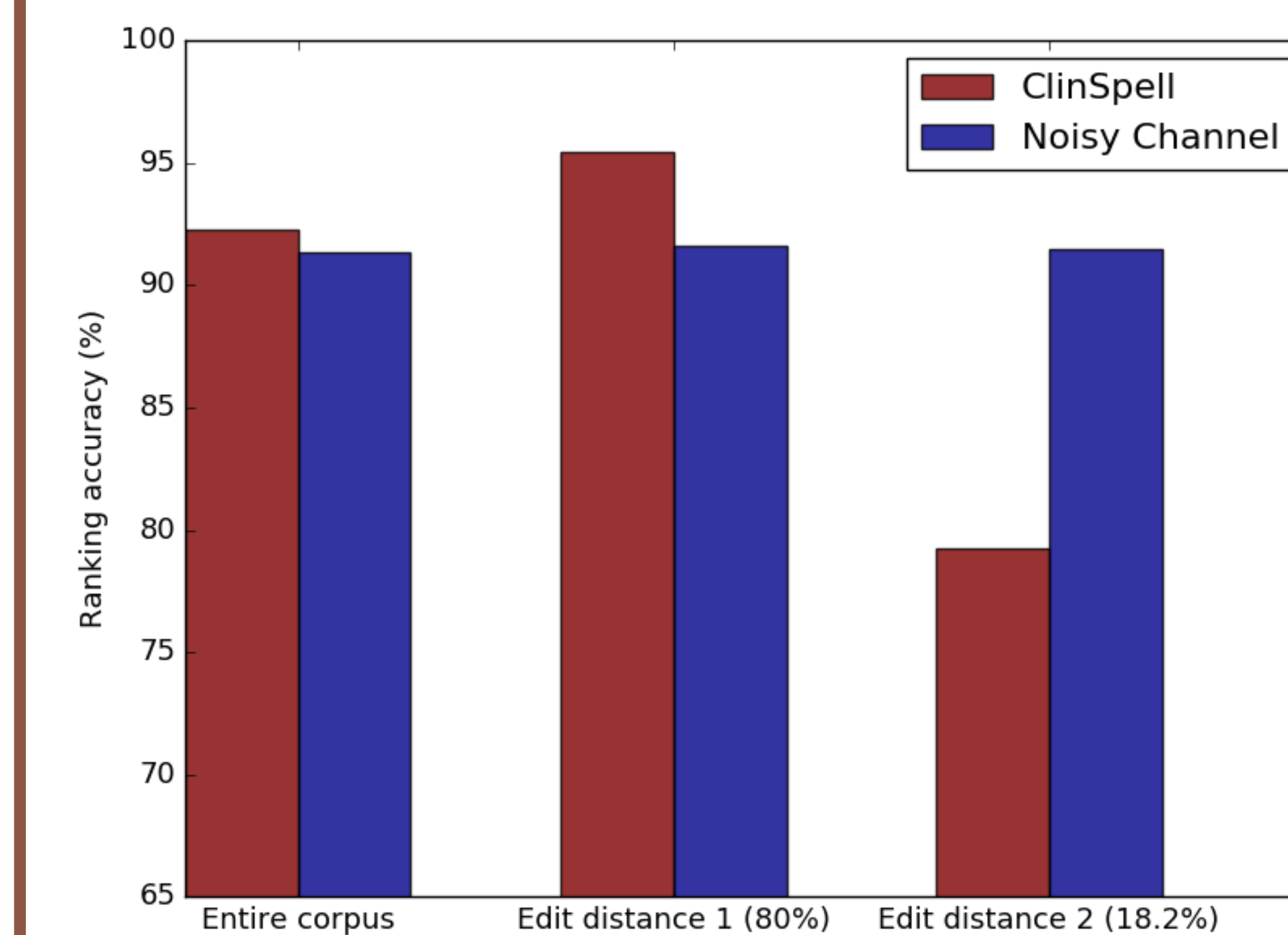


Figure 1: ENGLISH ranking accuracies

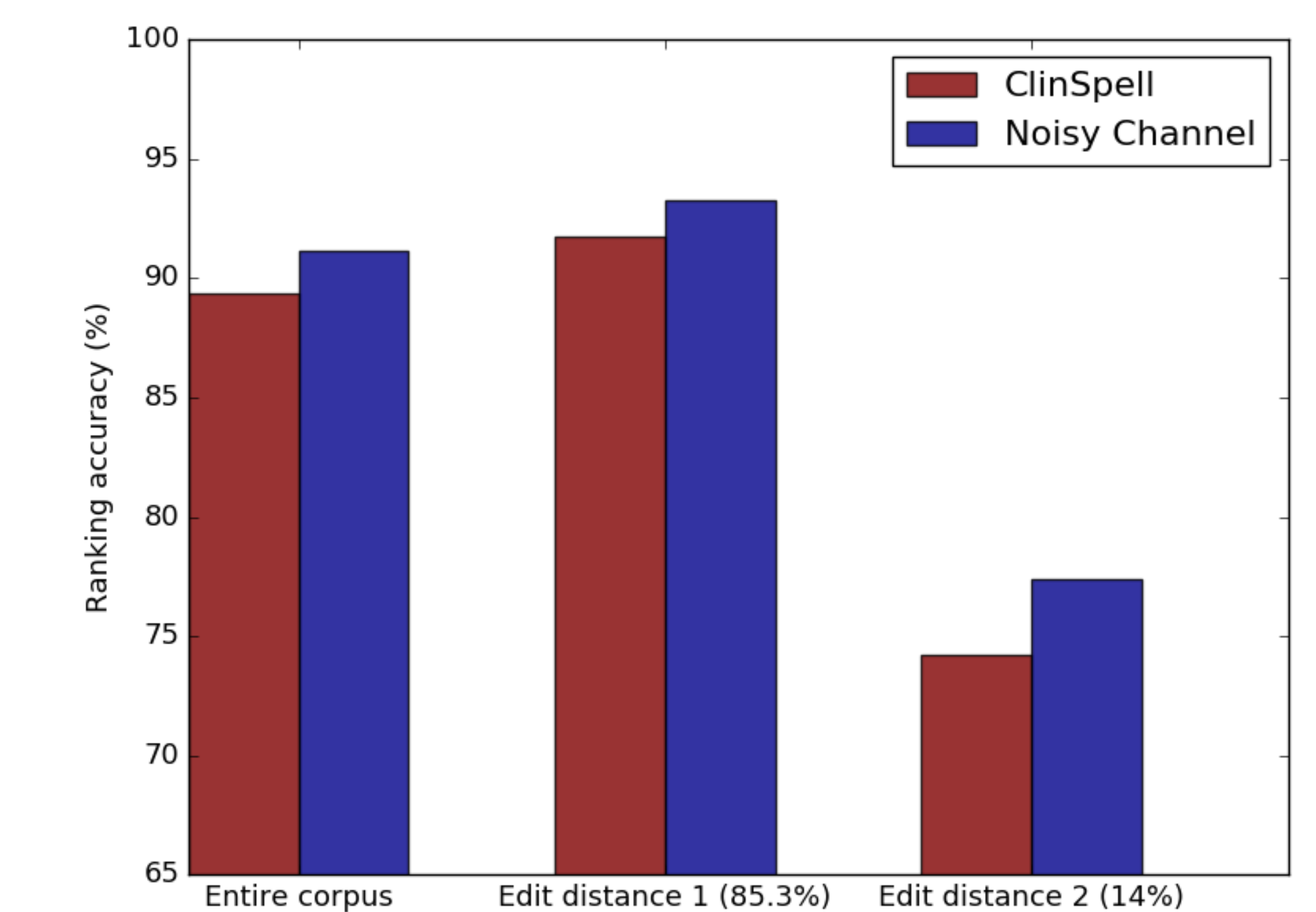


Figure 2: DUTCH ranking accuracies

CONCLUSIONS

- + Salient **performance gain** for English misspellings of 1 edit distance
- + **Contextual clues** can counter frequency bias

- Comparable performance when directly transferred to Dutch, but no superior performance yet → parameter tuning