



# Strong Baselines for Complex Word Identification across Multiple Languages



Pierre Finnimore, Elisabeth Fritsch, Daniel King, Alison Sneyd, Aneeq Ur Rehman, **Fernando Alva-Manchego**, Andreas Vlachos

## DATASETS (SECOND CWI SHARED TASK)

| Dataset / Language       | Train  | Dev   | Test  |
|--------------------------|--------|-------|-------|
| English (EN) - News      | 14,002 | 1,764 | 2,095 |
| English (EN) - WikiNews  | 7,746  | 870   | 1,287 |
| English (EN) - Wikipedia | 5,551  | 694   | 870   |
| Spanish (ES)             | 13,750 | 1,622 | 2,232 |
| German (DE)              | 6,151  | 795   | 959   |
| French (FR)              | N/A    | N/A   | 2,251 |

## MONOLINGUAL MODELS

**Ours:** Logistic Regression (LR) + 25 features

|                       |                                                                                                                                        |
|-----------------------|----------------------------------------------------------------------------------------------------------------------------------------|
| <b>Target-based</b>   | Named Entity type, part-of-speech, num. hypernym, num. tokens, language-normalised num. characters in each word, unigram probabilities |
| <b>Subword-based</b>  | prefixes, suffixes, num. syllables, num. complex punctuation marks (e.g. hyphens)                                                      |
| <b>Sentence-based</b> | sentence length, n-grams from the whole sentence                                                                                       |

| Dataset      | Ours        | Best CWI ST | SotA        |
|--------------|-------------|-------------|-------------|
| EN-News      | 86.0        | <b>87.4</b> | <b>87.4</b> |
| EN-WikiNews  | 81.6        | <b>84.0</b> | <b>84.0</b> |
| EN-Wikipedia | 76.1        | <b>81.2</b> | <b>81.2</b> |
| ES           | <b>77.6</b> | 77.0        | 77.0        |
| DE           | 74.8        | 74.5        | <b>75.5</b> |

same simple model

different complex models

Few language-independent features and Logistic Regression achieve comparable results to SotA complex models in monolingual and cross-lingual Complex Word Identification

| Sentence                                                                | Target Word or Multi-word Expression | Complex? |
|-------------------------------------------------------------------------|--------------------------------------|----------|
| <i>Both China and the Philippines flexed their muscles on Wednesday</i> | flexed                               | Yes      |
|                                                                         | flexed their muscles                 | Yes      |
|                                                                         | muscles                              | No       |

## Cross-lingual

### Features:

- number of syllables in the target
- number of tokens in the target
- number of complex punctuation marks (e.g. hyphens)
- sentence length
- unigram probabilities

| Dataset | Ours | SotA |
|---------|------|------|
| FR      | 75.8 | 76.0 |

LR + 5 Features      MTL + NNs + Random Forests Ensemble

## CROSS-LINGUAL MODELS

**Ours:** LR + 5 language-independent features

| EN | ES | DE | Test         | Macro F1 |
|----|----|----|--------------|----------|
|    |    | X  | EN-WikiNews  | 62.3     |
|    | X  |    | EN-Wikipedia | 63.1     |
|    | X  |    | EN-News      | 67.2     |
|    |    | X  | ES           | 72.6     |
| X  | X  |    | DE           | 73.4     |
|    |    | X  | FR           | 75.8     |

Adding EN to the training languages decreases performance for both ES and FR, but not for DE.

## DATASETS ANALYSIS

Every sub-word in 25% of multi-word expression instances has the opposite label

Non-complex — green  
 Complex — red

Every target multi-word expression in Spanish, German and French is labelled as complex



<https://github.com/sheffieldnlp/cwi>



@feralvam