# Multilingual Wikipedia:
# Editors of Primary Language Contribute to More Complex Articles

Sungjoon Park, Suin Kim*, Scott A. Hale, Sooyoung Kim, Jeongmin Byun, Alice Oh

{sungjoon.park, suin.kim, sooyoungkim, jmbyun}@kaist.ac.kr, scott.hale@oii.ox.ac.uk, alice.oh@kaist.edu

## Research Questions

- Is it possible to *quantify the complexity of language* in Wikipedia articles?
- Do primary users tend to edit parts of the Wikipedia articles with *higher language complexity* than non-primary users?
- Do we observe *more natural language* in the articles after primary users' edits compared to the articles after non-primary users' edits?
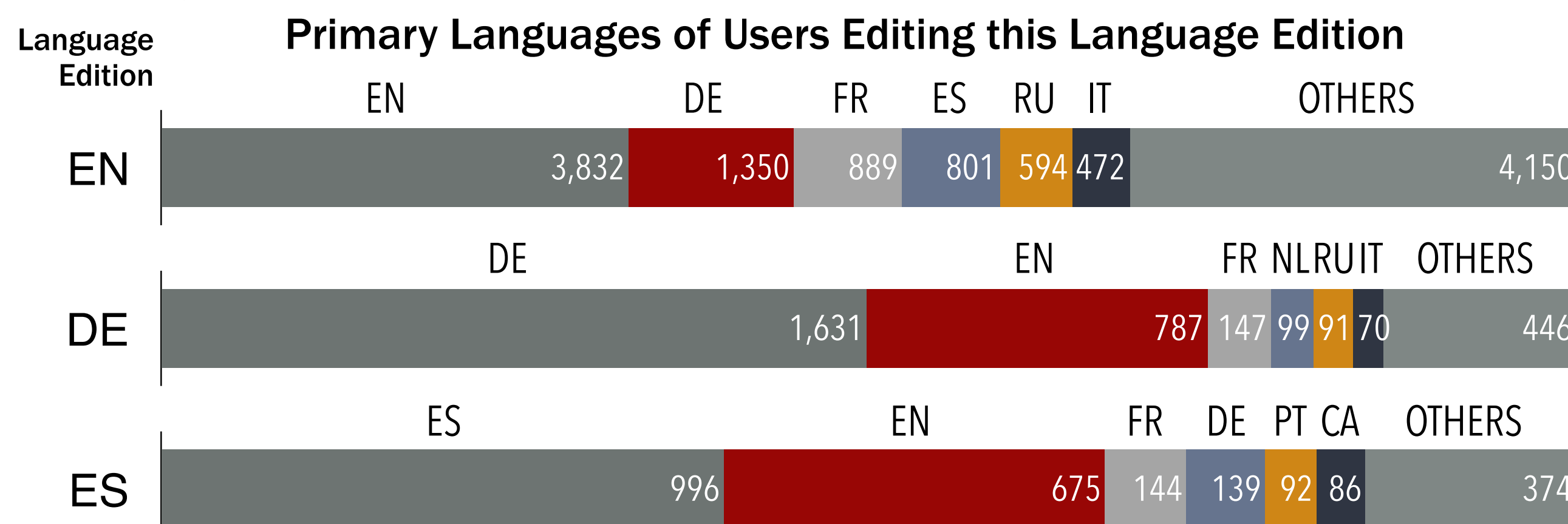
## Primary/Non-Primary Users

- Multilingual editors have edited different language editions
- We define a user's primary language as the language that the user edited the most (not necessarily user's native language)
- For each language X edition, we define primary users as those whose primary language is X
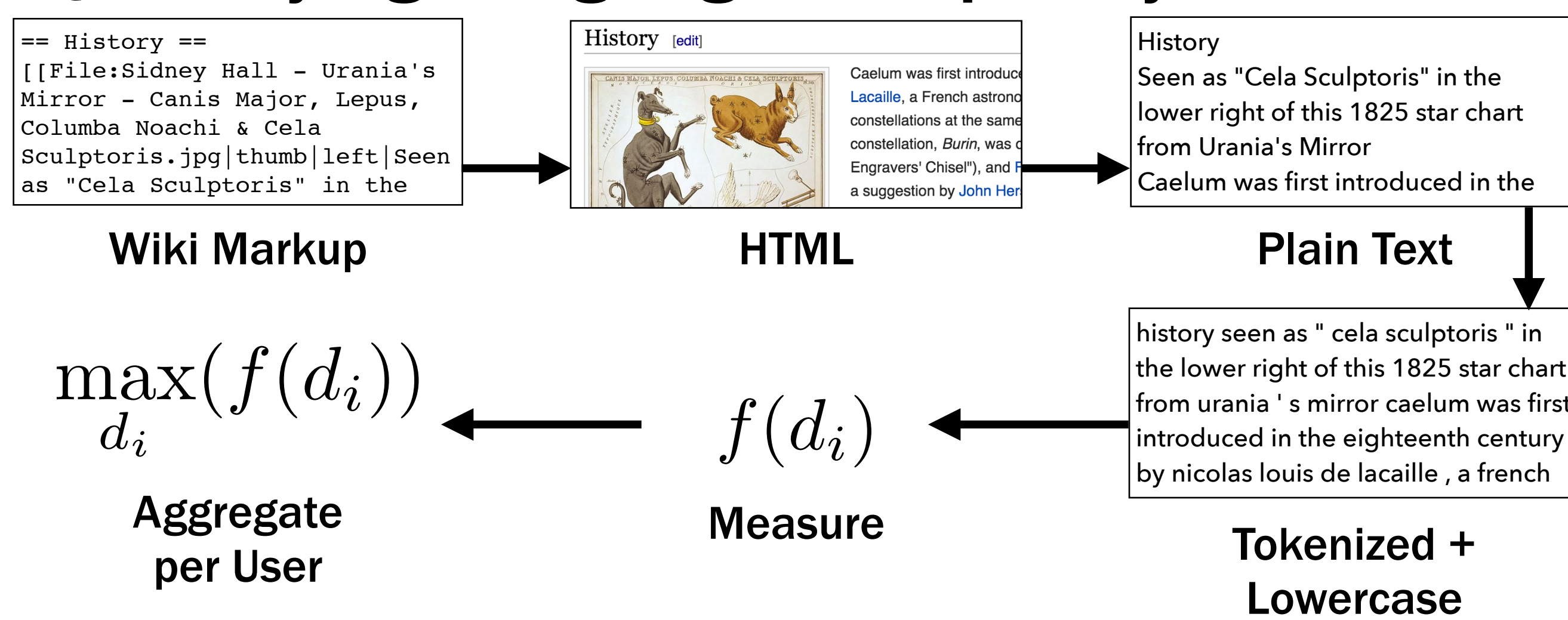
## Contributions

- We tested 20 measures for language complexity and showed they are highly consistent.
- Multilinguals in Wikipedia show relatively high levels of proficiency in their primary languages.

## Data

**Primary Languages of Users Editing this Language Edition**



| | English | German | Spanish |
|---|---|---|---|
| **#Editors** | 11,616 | 3,271 | 2,506 |
| **#Article Edit Sessions** | 237,849 | 120,123 | 69,557 |
| **#Edits** | 350,541 | 160,126 | 112,099 |

## Quantifying Language Complexity



**Wiki Markup** → **HTML** → **Plain Text**

$$f(d_i)$$
**Measure**

$$\max_{d_i}(f(d_i))$$
**Aggregate per User**

**Tokenized + Lowercase**

### Basic Features
- Number of characters
- Number of words
- Number of unique words
- Number of sentences
- Average word length in characters
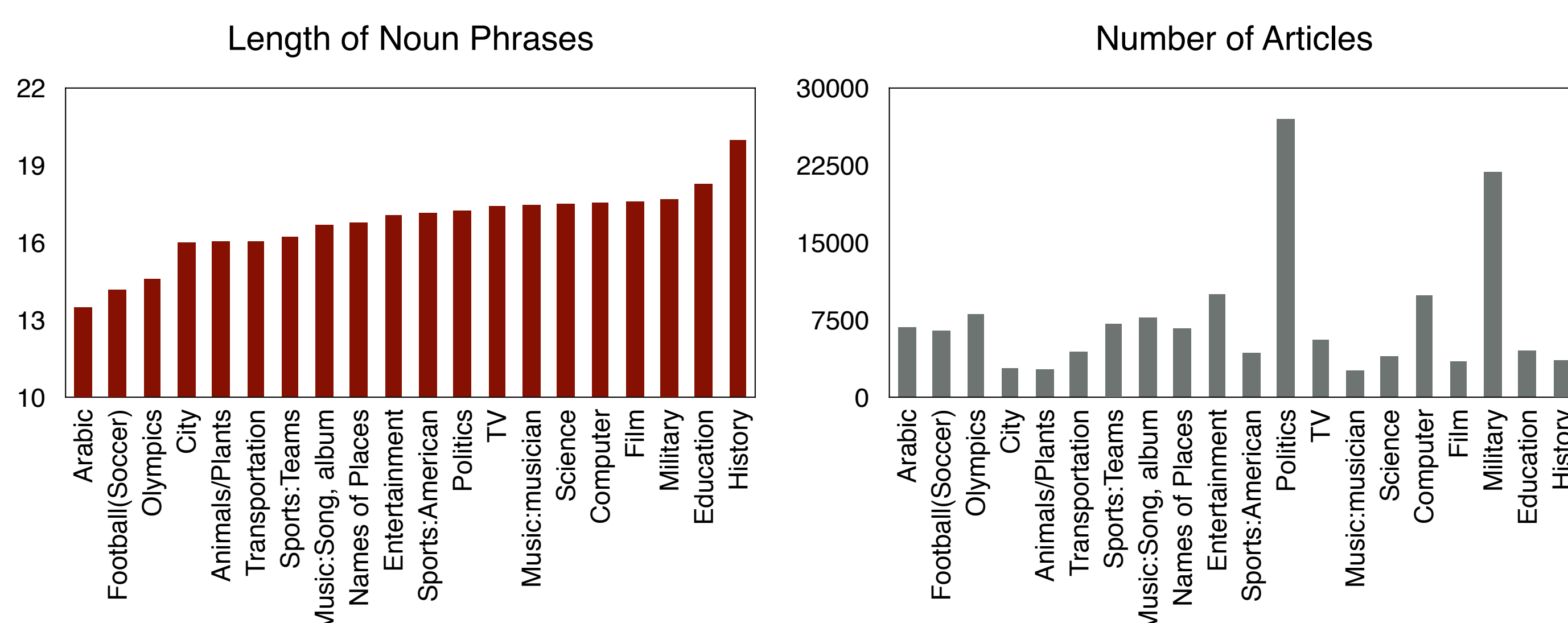- Average sentence length in words

### Lexical Diversity
- Entropy of word frequency
- Average word rank
- Average word occurrence frequency
- Error rate

### Syntactic Structure
- Entropy of POS frequency
- Mean phrase length (NP, VP)
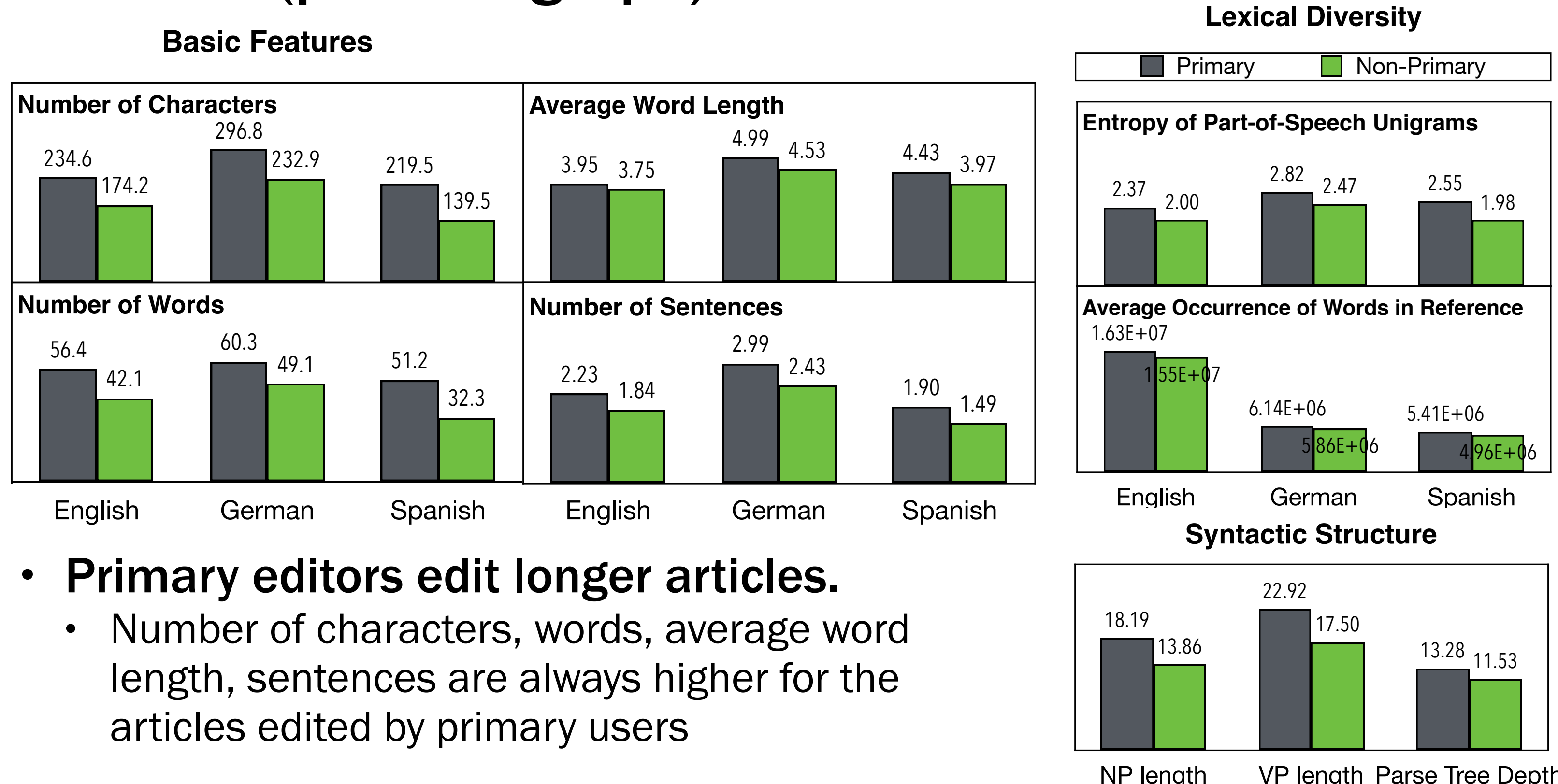- Mean parse tree depth

## Topic Clusters



Length of Noun Phrases

Number of Articles

- **LDA+DBSCAN Algorithm**
  - Fit LDA over entire data with 100 topics
    - Each document is reduced to 100 dimensions
  - Cluster documents into 20 clusters with DBSCAN
- **Cluster-level rank of language complexity is consistent over measures**
  - Articles in *history* cluster are more complex Articles in *football* cluster are less complex
  - The complexity measures used are consistent
  - Different topics show different complexity

## Results (per Paragraph)

**Basic Features**



Number of Characters

| English | German | Spanish |
|---|---|---|
| 234.6 / 174.2 | 296.8 / 232.9 | 219.5 / 139.5 |

Average Word Length

| English | German | Spanish |
|---|---|---|
| 3.95 / 3.75 | 4.99 / 4.53 | 4.43 / 3.97 |

Number of Words

| English | German | Spanish |
|---|---|---|
| 56.4 / 42.1 | 60.3 / 49.1 | 51.2 / 32.3 |

Number of Sentences

| English | German | Spanish |
|---|---|---|
| 2.23 / 1.84 | 2.99 / 2.43 | 1.90 / 1.49 |

**Lexical Diversity**

Primary / Non-Primary

Entropy of Part-of-Speech Unigrams

| English | German | Spanish |
|---|---|---|
| 2.37 / 2.00 | 2.82 / 2.47 | 2.55 / 1.98 |

Average Occurrence of Words in Reference

| English | German | Spanish |
|---|---|---|
| 1.63E+07 / 1.55E+07 | 6.14E+06 / 5.86E+06 | 5.41E+06 / 4.96E+06 |

Syntactic Structure

| NP length | VP length | Parse Tree Depth |
|---|---|---|
| 18.19 / 13.86 | 22.92 / 17.50 | 13.28 / 11.53 |

- **Primary editors edit articles with higher lexical diversity.**
  - Entropy of n-gram is always higher in primary edits
  - Primary users edit articles of frequent yet diverse set of words
- **Primary editors edit articles with more complex syntactic structure.**
  - Parse-tree based measures
    - Length of longest noun/verb phrase
    - Depth of parse tree
  - Articles edited by primary users have complex syntactic structure

- **Primary editors edit longer articles.**
  - Number of characters, words, average word length, sentences are always higher for the articles edited by primary users