CENTRE FOR DIVERSITY AND LEARNING, MULTIPLES, LANGUAGE & TRANSLATION TECHNOLOGY TEAM (LT3)

**Amaury Van Parys, Vanessa De Wilde, Lieve Macken, Maribel Montero Perez**

# LEXPRO: A PLURILINGUAL LEXICAL PROFILING TOOL FOR RESEARCH AND MATERIALS DEVELOPMENT

## Theoretical background

### Vocabulary
- Essential predictor of L2 reading comprehension (Jeon & Yamashita, 2022)
- More words known = better comprehension (Schmitt et al., 2011)
- Need for objective method to assess how demanding a text's vocabulary will be for L2 learners

### Lexical profiling
- Method for determining **vocabulary demands** of L2 input
- Often used in previous research: TV series (Webb & Rodgers, 2009), novels (Nation, 2006), L2 textbooks (Van Parys et al., 2024), etc.
- Categorising vocabulary across **word frequency** levels: Higher-frequency words have higher odds of being known by learner (Nation, 2013)
- Allows to estimate **vocabulary loads**, i.e., estimates of required vocabulary size for achieving crucial points of vocabulary coverage (Webb, 2020):
  - **95%** coverage: needed for basic comprehension
  - **98%** coverage: needed for detailed comprehension
- Example: according to Webb & Rodgers (2009), the 3,000 most frequent word families in English need to be known for 95% coverage of TV series and thus basic comprehension

## Gaps in prior profiling methods

### Word counting unit
- Typical counting unit in profiling is the **word family**, which covers a headword (e.g., 'act') with all its inflections (e.g., 'acts', 'acting') and derivations (e.g., 'actor')
- However: increasing criticism (e.g., McLean, 2018; Stoeckel et al., 2024)
- Potentially more appropriate counting unit: the **flemma,** which covers a headword with all inflections (across different parts of speech), but not derivations

### Word frequency as proxy for learner knowledge
- Typically used frequency lists in profiling are based on **broad corpora** covering a wide range of topics (e.g., British National Corpus)
- However: these lists do not reflect learner knowledge as closely as once presumed (Pinchbeck et al., 2022)
- Lists derived from **subtitle corpora** (e.g., SubtLex-UK) appear to align more with learner knowledge (Pinchbeck et al., 2022; van Heuven et al., 2014)

### Lack of focus on non-English L2s
- Most existing tools (e.g., LexTutor; Cobb, n.d.) mainly target English
- In line with overall focus on English in Second Language Acquisition (Brezina & Pallotti, 2019)

## Goals of LexPro

LexPro aims to set itself apart from existing profiling tools by:
- Using the **flemma** as main word counting unit
- Using **subtitle-based frequency lists**
- Targeting **English** in addition to multiple other L2s (currently **French**, **Spanish**, and **Dutch**)
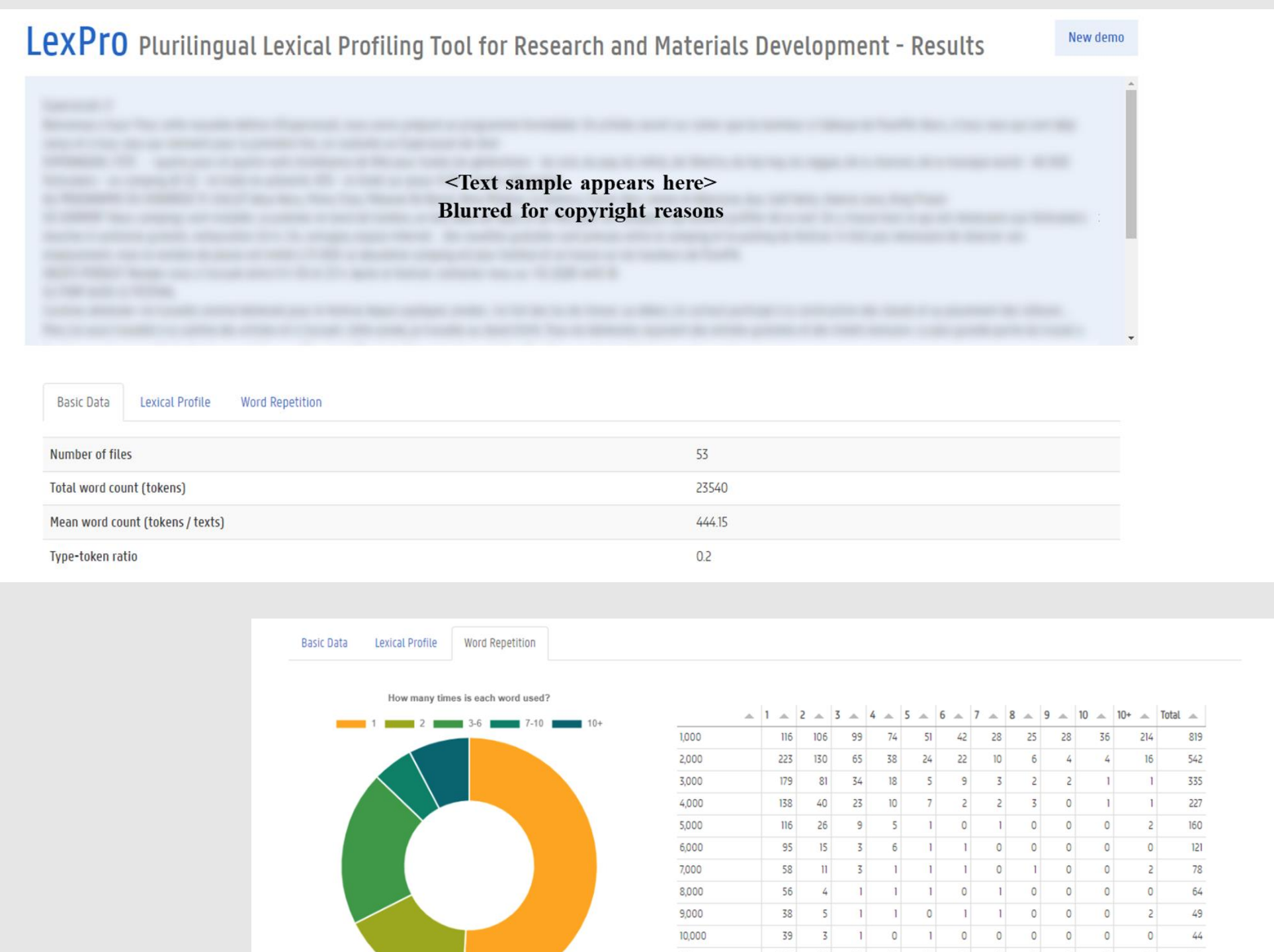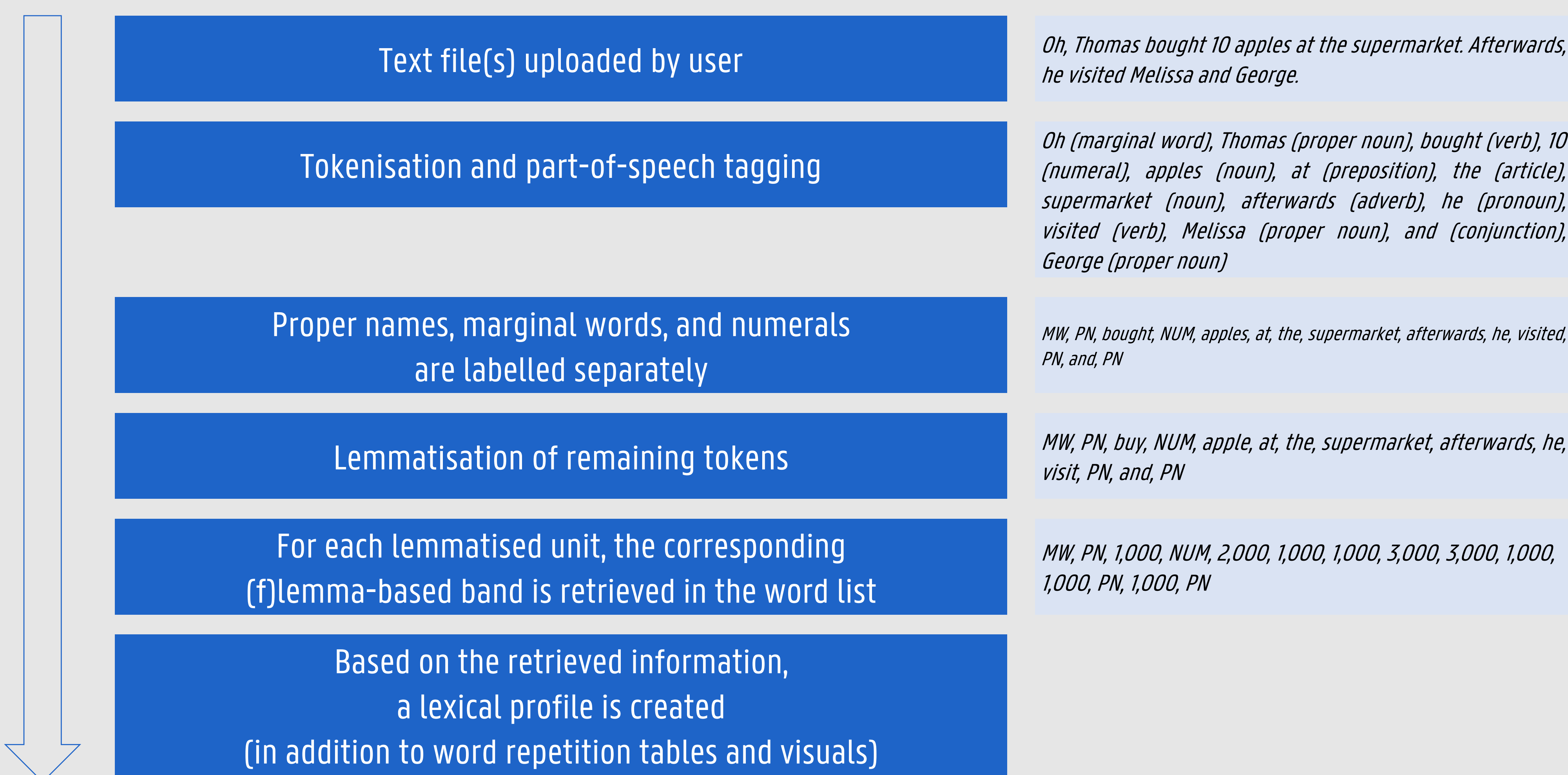
## Output

Example corpus: French Netflix series 'Lupin' (Kay & Uzan, 2021-present)
- 3 seasons
- 17 episodes
- 54,502 running tokens
- On average 3,206 tokens per episode

## How does LexPro work?

Programmed in **Python** and relying on the **spaCy** NLP library (Honnibal & Montani, 2017)

| Step | Example |
|---|---|
| Text file(s) uploaded by user | *Oh, Thomas bought 10 apples at the supermarket. Afterwards, he visited Melissa and George.* |
| Tokenisation and part-of-speech tagging | *Oh (marginal word), Thomas (proper noun), bought (verb), 10 (numeral), apples (noun), at (preposition), the (article), supermarket (noun), afterwards (adverb), he (pronoun), visited (verb), Melissa (proper noun), and (conjunction), George (proper noun)* |
| Proper names, marginal words, and numerals are labelled separately | *MW, PN, bought, NUM, apples, at, the, supermarket, afterwards, he, visited, PN, and, PN* |
| Lemmatisation of remaining tokens | *MW, PN, buy, NUM, apple, at, the, supermarket, afterwards, he, visit, PN, and, PN* |
| For each lemmatised unit, the corresponding (f)lemma-based band is retrieved in the word list | *MW, PN, 1,000, NUM, 2,000, 1,000, 1,000, 3,000, 3,000, 1,000, 1,000, PN, 1,000, PN* |
| Based on the retrieved information, a lexical profile is created (in addition to word repetition tables and visuals) | |

**GHENT UNIVERSITY**