# CYP-LEX:
# A novel large-scale lexical database
# of books for children and young people

**Maria Korochkina** | **Marco Marelli** | **Marc Brysbaert** | **Kathy Rastle**

*Royal Holloway University of London*

May 31, 2023

- Language is a *productive* system in which a *finite* set of elements can be combined to express a *limitless* range of thought

- Language is a *productive* system in which a *finite* set of elements can be combined to express a *limitless* range of thought

- Most words are built by combining units of meaning (morphemes) which contribute *systematically* to the meaning of the word

- Language is a *productive* system in which a *finite* set of elements can be combined to express a *limitless* range of thought

- Most words are built by combining units of meaning (morphemes) which contribute *systematically* to the meaning of the word
  - 'un-' + 'clean' → 'unclean' (not clean)

- Language is a *productive* system in which a *finite* set of elements can be combined to express a *limitless* range of thought

- Most words are built by combining units of meaning (morphemes) which contribute *systematically* to the meaning of the word
  - 'un-' + 'clean' → 'unclean' (not clean)
  - [stem] + '-ify' → 'make [stem]' (e.g., simplify, purify)

- Language is a *productive* system in which a *finite* set of elements can be combined to express a *limitless* range of thought

- Most words are built by combining units of meaning (morphemes) which contribute *systematically* to the meaning of the word
    - 'un-' + 'clean' → 'unclean' (not clean)
    - [stem] + '-ify' → 'make [stem]' (e.g., simplify, purify)

- This combinatorial system allows us to *generalise* and provides a heuristic for vocabulary growth

- Language is a *productive* system in which a *finite* set of elements can be combined to express a *limitless* range of thought

- Most words are built by combining units of meaning (morphemes) which contribute *systematically* to the meaning of the word
  - 'un-' + 'clean' → 'unclean' (not clean)
  - [stem] + '-ify' → 'make [stem]' (e.g., simplify, purify)

- This combinatorial system allows us to *generalise* and provides a heuristic for vocabulary growth
  - 'quick' + '-ify' → 'quickify'

- Affixes alter words' meanings in highly predictable ways [1, 2]
  & do not occur in isolation

- Affixes alter words' meanings in highly predictable ways [1, 2] & do not occur in isolation
- Their function can & must be inferred from experience with whole words (e.g., solidify, amplify, justify)

- Affixes alter words' meanings in highly predictable ways [1, 2] & do not occur in isolation

- Their function can & must be inferred from experience with whole words (e.g., solidify, amplify, justify)

- However, morphemic regularities are not as regular as one might think...

- Affixes alter words' meanings in highly predictable ways [1, 2] & do not occur in isolation

- Their function can & must be inferred from experience with whole words (e.g., solidify, amplify, justify)

- However, morphemic regularities are not as regular as one might think...

- There is substantial variability in precision & consistency with which a morpheme conveys meaningful information

- Affixes alter words' meanings in highly predictable ways [1, 2] & do not occur in isolation
- Their function can & must be inferred from experience with whole words (e.g., solidify, amplify, justify)

- However, morphemic regularities are not as regular as one might think...
- There is substantial variability in precision & consistency with which a morpheme conveys meaningful information
  - 'artist', 'typist' vs. 'racist', 'sadist'

- Affixes alter words' meanings in highly predictable ways [1, 2] & do not occur in isolation

- Their function can & must be inferred from experience with whole words (e.g., solidify, amplify, justify)

- However, morphemic regularities are not as regular as one might think...

- There is substantial variability in precision & consistency with which a morpheme conveys meaningful information
    - 'artist', 'typist' vs. 'racist', 'sadist'

→ **What** do we need to learn and **how**?

- Most new words children aged 8+ acquire are encountered through reading [3]

- Most new words children aged 8+ acquire are encountered through reading [3]
- In adults, morphological knowledge is modulated by reading experience [4]

- Most new words children aged 8+ acquire are encountered through reading [3]
- In adults, morphological knowledge is modulated by reading experience [4]
- In English, morphological information is far more salient in spelling than it is in spoken language (e.g., 'sign' vs. 'signature'; 'active' vs. 'activity') [5]

- Most new words children aged 8+ acquire are encountered through reading [3]
- In adults, morphological knowledge is modulated by reading experience [4]
- In English, morphological information is far more salient in spelling than it is in spoken language (e.g., 'sign' vs. 'signature'; 'active' vs. 'activity') [5]
- Yet, available corpora are too small, target narrow age cohorts, comprise outdated reading materials, or have restricted access

- Most new words children aged 8+ acquire are encountered through reading [3]
- In adults, morphological knowledge is modulated by reading experience [4]
- In English, morphological information is far more salient in spelling than it is in spoken language (e.g., 'sign' vs. 'signature'; 'active' vs. 'activity') [5]
- Yet, available corpora are too small, target narrow age cohorts, comprise outdated reading materials, or have restricted access

↓

We need a large-scale publicly available corpus of books that children and young people read!

National reading surveys, publisher data, & book sales statistics
from Amazon, BookTrust, Goodreads, LoveReading4Kids, etc.

National reading surveys, publisher data, & book sales statistics
from Amazon, BookTrust, Goodreads, LoveReading4Kids, etc.

⇓

1,200 popular fiction & non-fiction e-books

400 books per age band (7-9, 10-12, 13+)

National reading surveys, publisher data, & book sales statistics
from Amazon, BookTrust, Goodreads, LoveReading4Kids, etc.
⇓
1,200 popular fiction & non-fiction e-books
400 books per age band (7-9, 10-12, 13+)
⇓
Semi-automatic conversion, cleaning, & pre-processing

National reading surveys, publisher data, & book sales statistics
from Amazon, BookTrust, Goodreads, LoveReading4Kids, etc.

⇓

1,200 popular fiction & non-fiction e-books
400 books per age band (7-9, 10-12, 13+)

⇓

Semi-automatic conversion, cleaning, & pre-processing

⇓

Tokenisation, lemmatisation, Part-of-Speech tagging
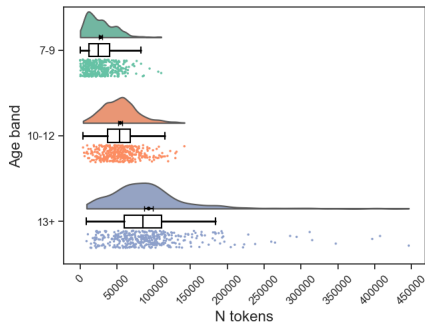71,296,249 tokens & 323,670 types

National reading surveys, publisher data, & book sales statistics
from Amazon, BookTrust, Goodreads, LoveReading4Kids, etc.
⇓
1,200 popular fiction & non-fiction e-books
400 books per age band (7-9, 10-12, 13+)
⇓
Semi-automatic conversion, cleaning, & pre-processing
⇓
Tokenisation, lemmatisation, Part-of-Speech tagging
71,296,249 tokens & 323,670 types
⇓
Invalid token removal

National reading surveys, publisher data, & book sales statistics
from Amazon, BookTrust, Goodreads, LoveReading4Kids, etc.
⇓
1,200 popular fiction & non-fiction e-books
400 books per age band (7-9, 10-12, 13+)
⇓
Semi-automatic conversion, cleaning, & pre-processing
⇓
Tokenisation, lemmatisation, Part-of-Speech tagging
71,296,249 tokens & 323,670 types
⇓
Invalid token removal
⇓
**CYP-LEX: Children and Young People's Books Lexicon**
70,287,217 tokens & 105,694 types

| | CYP-LEX 7–9 | CYP-LEX 10–12 | CYP-LEX 13+ |
|---|---|---|---|
| $N$ tokens | 11,162,653 | 21,837,794 | 37,286,770 |
| $\mu$ $(\sigma)$ tokens | 27,906.63 (19,212.43) | 54,594.48 (24,011.91) | 93,216.92 (57,718.38) |
| $N$ types | 52,851 | 70,945 | 90,980 |
| $\mu$ $(\sigma)$ types | 3,027.81 (1,452.05) | 4,712.81 (1,550.43) | 6,446.57 (2,365.59) |

| | Cbeebies 0–6 years $N = 27{,}236$ | | CPWD 5–9 years $N = 12{,}452$ | | CBBC 6–12 years $N = 58{,}691$ | | SUBTLEX-UK Adults $N = 160{,}024$ | |
|---|---|---|---|---|---|---|---|---|
| | % shared | r | % shared | r | % shared | r | % shared | r |
| **CYP-LEX 7–9** $N = 52{,}851$ | 39% | .67 | 19% | .71 | 70% | .77 | 91% | .72 |
| **CYP-LEX 10–12** $N = 70{,}945$ | 30% | .63 | 14% | .68 | 58% | .75 | 86% | .76 |
| **CYP-LEX 13+** $N = 90{,}980$ | 24% | .58 | 11% | .62 | 48% | .72 | 79% | .76 |

| | Cbeebies 0–6 years N = 27,236 | | CPWD 5–9 years N = 12,452 | | CBBC 6–12 years N = 58,691 | | SUBTLEX-UK Adults N = 160,024 | |
|---|---|---|---|---|---|---|---|---|
| | % shared | r | % shared | r | % shared | r | % shared | r |
| **CYP-LEX 7–9** N = 52,851 | 39% | .67 | 19% | .71 | 70% | .77 | 91% | .72 |
| **CYP-LEX 10–12** N = 70,945 | 30% | .63 | 14% | .68 | 58% | .75 | 86% | .76 |
| **CYP-LEX 13+** N = 90,980 | 24% | .58 | 11% | .62 | 48% | .72 | 79% | .76 |

|  | Cbeebies 0–6 years $N = 27{,}236$ | | CPWD 5–9 years $N = 12{,}452$ | | CBBC 6–12 years $N = 58{,}691$ | | SUBTLEX-UK Adults $N = 160{,}024$ | |
|---|---|---|---|---|---|---|---|---|
|  | % shared | $r$ | % shared | $r$ | % shared | $r$ | % shared | $r$ |
| **CYP-LEX 7–9** $N = 52{,}851$ | 39% | .67 | 19% | .71 | 70% | .77 | 91% | .72 |
| **CYP-LEX 10–12** $N = 70{,}945$ | 30% | .63 | 14% | .68 | 58% | .75 | 86% | .76 |
| **CYP-LEX 13+** $N = 90{,}980$ | 24% | .58 | 11% | .62 | 48% | .72 | 79% | .76 |

|  | Cbeebies 0–6 years $N = 27{,}236$ | | CPWD 5–9 years $N = 12{,}452$ | | CBBC 6–12 years $N = 58{,}691$ | | SUBTLEX-UK Adults $N = 160{,}024$ | |
|---|---|---|---|---|---|---|---|---|
|  | % shared | $r$ | % shared | $r$ | % shared | $r$ | % shared | $r$ |
| **CYP-LEX 7–9** $N = 52{,}851$ | 39% | .67 | 19% | .71 | 70% | .77 | 91% | .72 |
| **CYP-LEX 10–12** $N = 70{,}945$ | 30% | .63 | 14% | .68 | 58% | .75 | 86% | .76 |
| **CYP-LEX 13+** $N = 90{,}980$ | 24% | .58 | 11% | .62 | 48% | .72 | 79% | .76 |

# CYP-LEX 10–12 vs. CYP-LEX 7–9

25,627 unshared words

# CYP-LEX 10–12 vs. CYP-LEX 7–9

25,627 unshared words

Raw frequency $\leq 3$
73% ($N = 18{,}646$) of unshared words

# CYP-LEX 10–12 vs. CYP-LEX 7–9

25,627 unshared words

Raw frequency $\leq 3$
73% ($N = 18,646$) of unshared words

CYP-LEX 10–12 vs. CYP-LEX 7–9

25,627 unshared words

Raw frequency ≤ 3
73% ($N = 18,646$) of unshared words

Raw frequency > 100
< 1% ($N = 249$) of unshared words

CYP-LEX 10–12 vs. CYP-LEX 7–9

25,627 unshared words

Raw frequency ≤ 3
73% ($N = 18,646$) of unshared words

Raw frequency > 100
< 1% ($N = 249$) of unshared words

# CYP-LEX 13+ vs. CYP-LEX 10–12

31,025 unshared words

CYP-LEX 13+ vs. CYP-LEX 10–12

31,025 unshared words

Raw frequency $\leq 3$
74% ($N = 22,855$) of unshared words

CYP-LEX 13+ vs. CYP-LEX 10–12

31,025 unshared words

Raw frequency $\leq 3$
74% ($N = 22{,}855$) of unshared words

Raw frequency $> 100$
1% ($N = 326$) of unshared words

# CYP-LEX 13+ vs. CYP-LEX 10–12

31,025 unshared words

Raw frequency $\leq 3$
74% ($N = 22{,}855$) of unshared words

Raw frequency $> 100$
1% ($N = 326$) of unshared words

CYP-LEX 13+ vs. CYP-LEX 10–12

31,025 unshared words

Raw frequency $\leq 3$
74% ($N = 22{,}855$) of unshared words

Raw frequency $> 100$
1% ($N = 326$) of unshared words



On that note...

CYP-LEX 13+ vs. CYP-LEX 10–12

31,025 unshared words

Raw frequency $\leq 3$
74% ($N = 22{,}855$) of unshared words

Raw frequency $> 100$
1% ($N = 326$) of unshared words





On that note...

- CYP-LEX 13+ vs. CBBC: 'shit' (*Zipf* = 4.92, $N = 3{,}077$)

# CYP-LEX 13+ vs. CYP-LEX 10–12

31,025 unshared words

Raw frequency $\leq 3$
74% ($N = 22{,}855$) of unshared words

Raw frequency $> 100$
1% ($N = 326$) of unshared words





On that note...

- CYP-LEX 13+ vs. CBBC: 'shit' ($Zipf = 4.92$, $N = 3{,}077$)
- CYP-LEX 13+ vs. CPWD: 'hell' ($Zipf = 5.21$, $N = 6{,}103$)

7–9    10–12    13+

# Semantic similarity across the age bands

### 600 most common words in sets of 100

All lemmas — Excluding function words

↓

The 'older' the children, the more similar the books?..

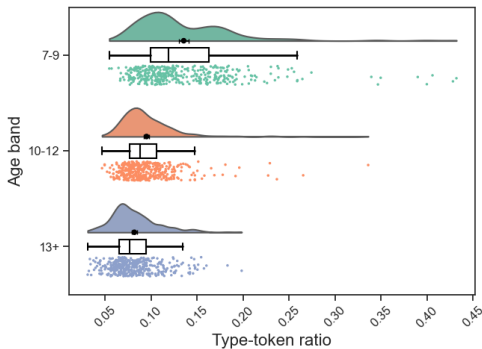- Large corpora typically have higher TTR: the longer the sample, the higher the probability of encountering a new word

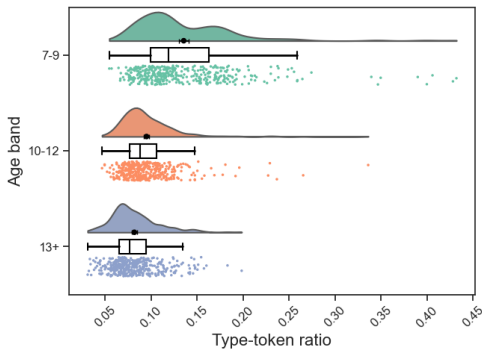- Large corpora typically have higher TTR: the longer the sample, the higher the probability of encountering a new word

- Large corpora typically have higher TTR: the longer the sample, the higher the probability of encountering a new word
- In CYP-LEX, this relationship is inverse!

- 100 most frequent vs. infrequent lemmas across all the books in the 7–9 age band

- 100 most frequent vs. infrequent lemmas across all the books in the 7–9 age band

- 100 most frequent vs. infrequent lemmas across all the books in the 7–9 age band

Most of these occur in all books

- 100 most frequent vs. infrequent lemmas across all the books in the 7–9 age band

Most of these occur in all books

- 100 most frequent vs. infrequent lemmas across all the books
  in the 7–9 age band

Most of these occur in all books

These ones occur in 3 books only

*Thank you!*

[1]  D. Plaut and L. Gonnerman. "Are non-semantic
     morphological effects incompatible with a distributed
     connectionist approach to lexical processing?" In: *Language
     and Cognitive Processes* 15 (2000), pp. 445–485. DOI:
     https://doi.org/10.1080/01690960050119661.

[2]  K. Rastle and M. H. Davis. "Morphological decomposition
     based on the analysis of orthography". In: *Language and
     Cognitive Processes* 23 (2008), pp. 942–971. DOI:
     https://doi.org/10.1080/01690960802069730.

[3]  W. Nagy and R. Andreson. "How many words are there in
     printed school English?" In: *Reading Research Quarterly* 19
     (1984), pp. 304–330. DOI: https://doi.org/10.2307/747823.

[4]  S. Andrews and S. Lo. "Is morphological priming stronger for transparent than opaque words? It depends on individual differences in spelling and vocabulary". In: *Journal of Memory and Language* 68 (2013), pp. 279–296. DOI: https://doi.org/10.1016/j.jml.2012.12.001.

[5]  A. Ulicheva et al. "Skilled readers' sensitivity to meaningful regularities in English writing". In: *Cognition* 195 (2020), p. 103810. DOI: https://doi.org/10.1016/j.cognition.2018.09.013.