# Corpus Research on the Development of Children's School Writing

Phil Durrant

**University of Exeter**

UNIVERSITY OF EXETER

E·S·R·C
ECONOMIC
& SOCIAL
RESEARCH
COUNCIL

That morning a meerkat mob was snoring there heds of Suddenly a snaka slithed into the brow. The snaka saw the baby. The baby ran to the mum, the mum froo the snak. Next a jackal ran in to Sunny. The jackal sede can I be your frend? And they play together tag. Sunny a vitid him for tea time the end.

one luge time ago there was a king colld king james the first and the cathlixs did not like him. and there was a bad man called Guy Fawkes he wantied to bow the houses of Parliament he wantid to cill the king to as well as the cathlixs he had 36 barols of gunpowder and he hid it. Robert Catesby sent a leter to the king.
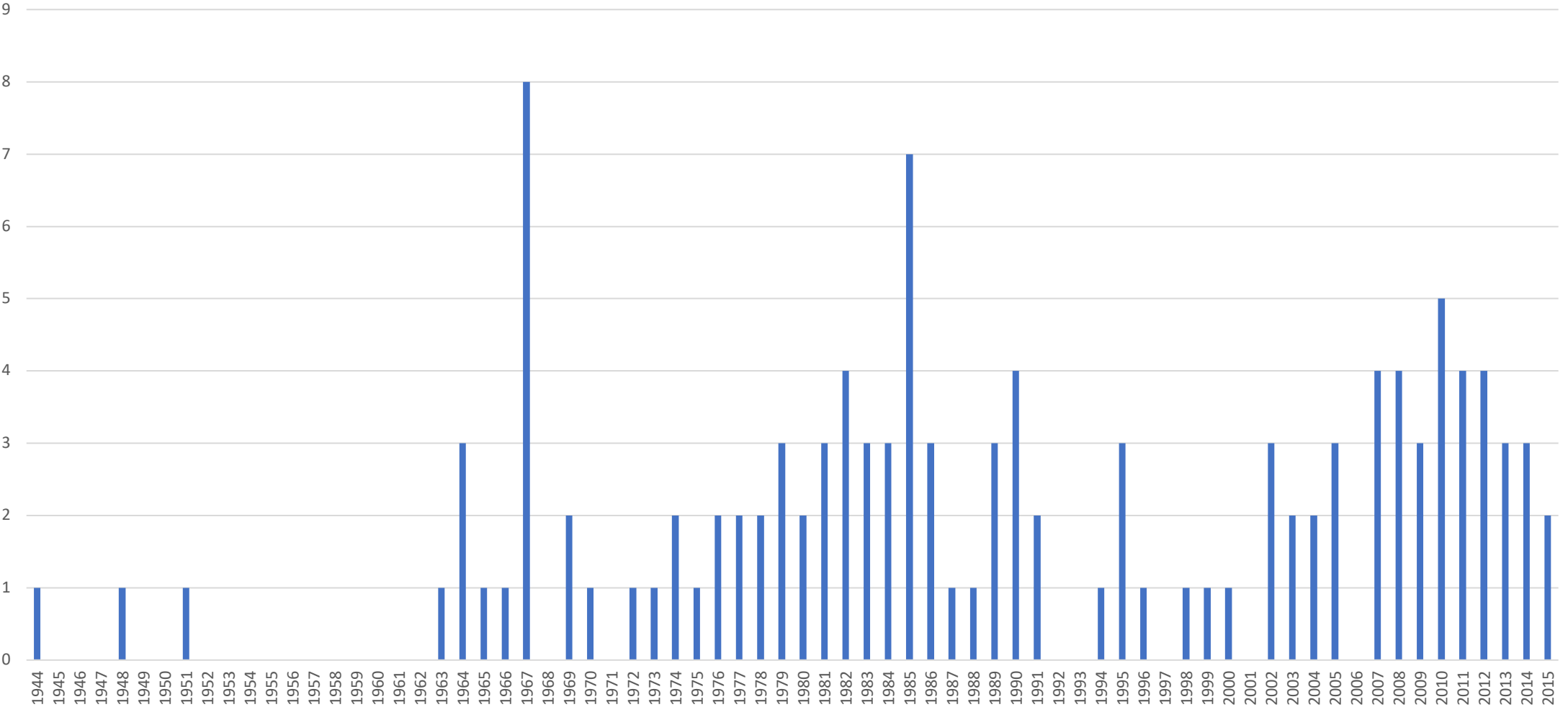
Dear Sir, I am writing to express my views on the article you recently printed, detailing a scheme by the Divert Trust to help difficult students. At first I was unsure if this scheme could ever work, and was indignant, like so many others, that many good students remained unrewarded. However, after researching this scheme I have come to realise that it is rather a brilliant idea. Research shows that around 88% of schools admit to not being able to cope with difficult students.

# Historical Overview

Quantitative comparisons of linguistic features in mainstream children's writing
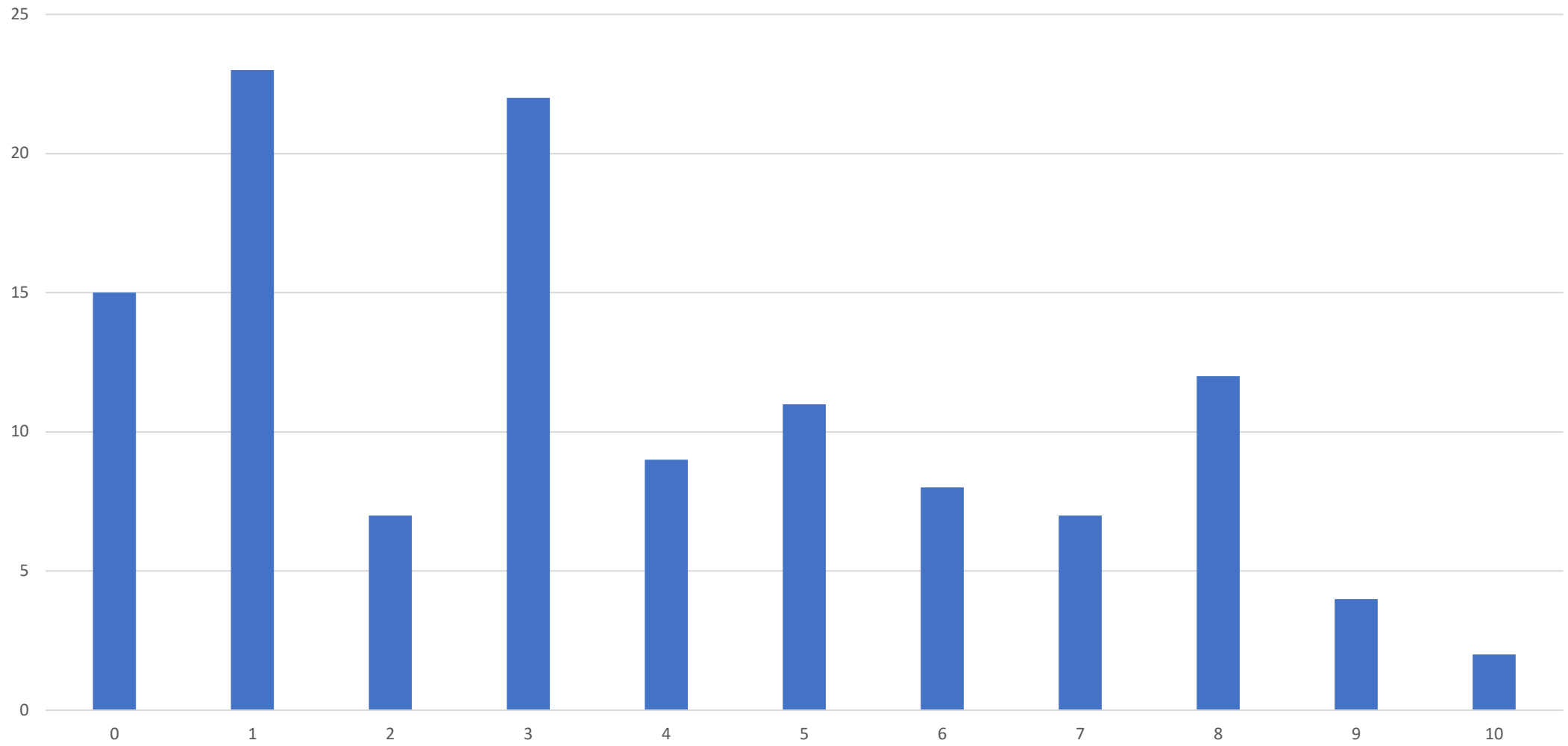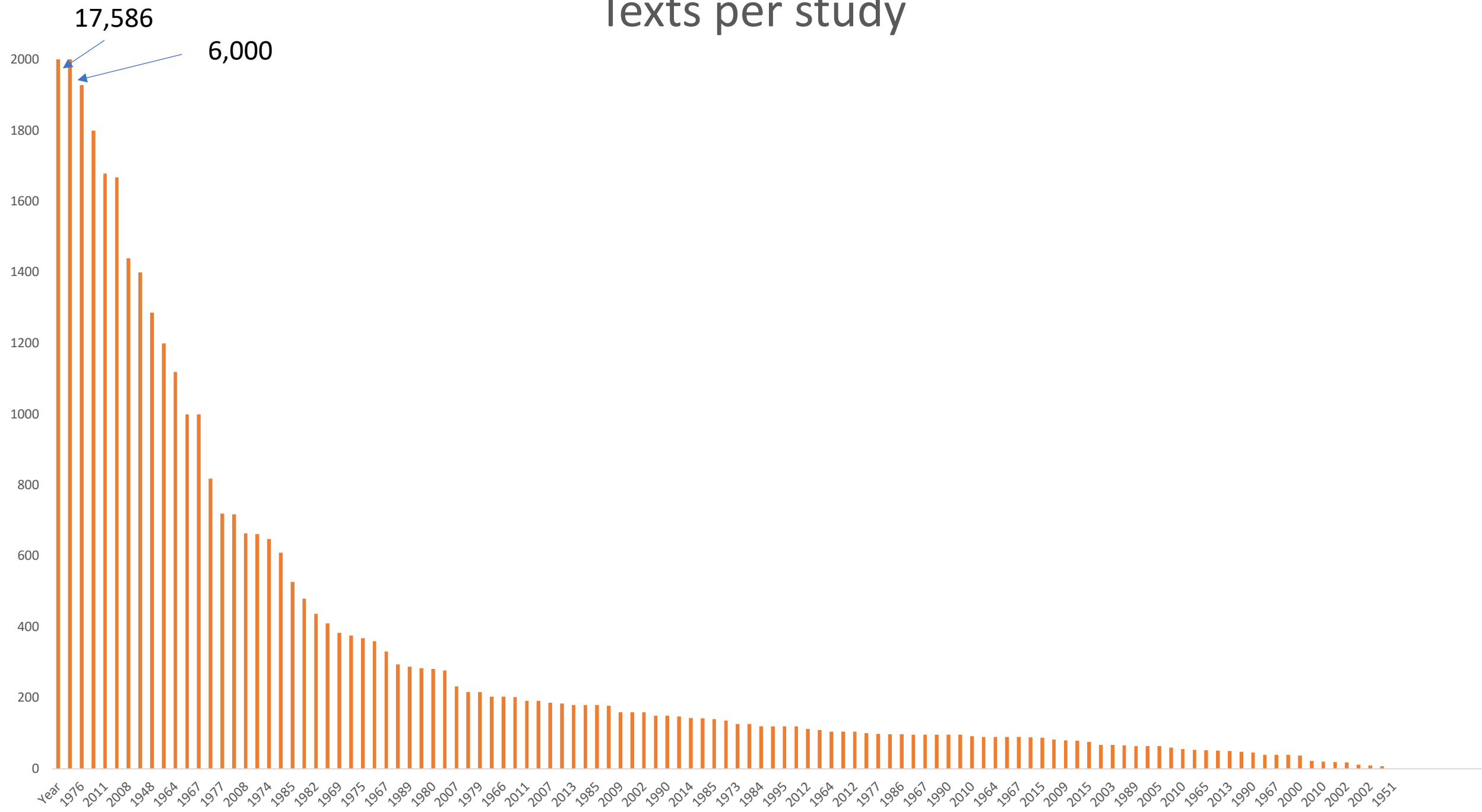
Studies per year: 1944-2015
(Total 120 studies)

# Provenance

| Country | Number of Studies |
|---|---:|
| USA | 77 |
| UK | 22 |
| Canada | 16 |
| Australia | 2 |
| USA, UK & New Zealand | 2 |
| USA & New Zealand | 1 |

Age range

Texts per study

17,586

6,000
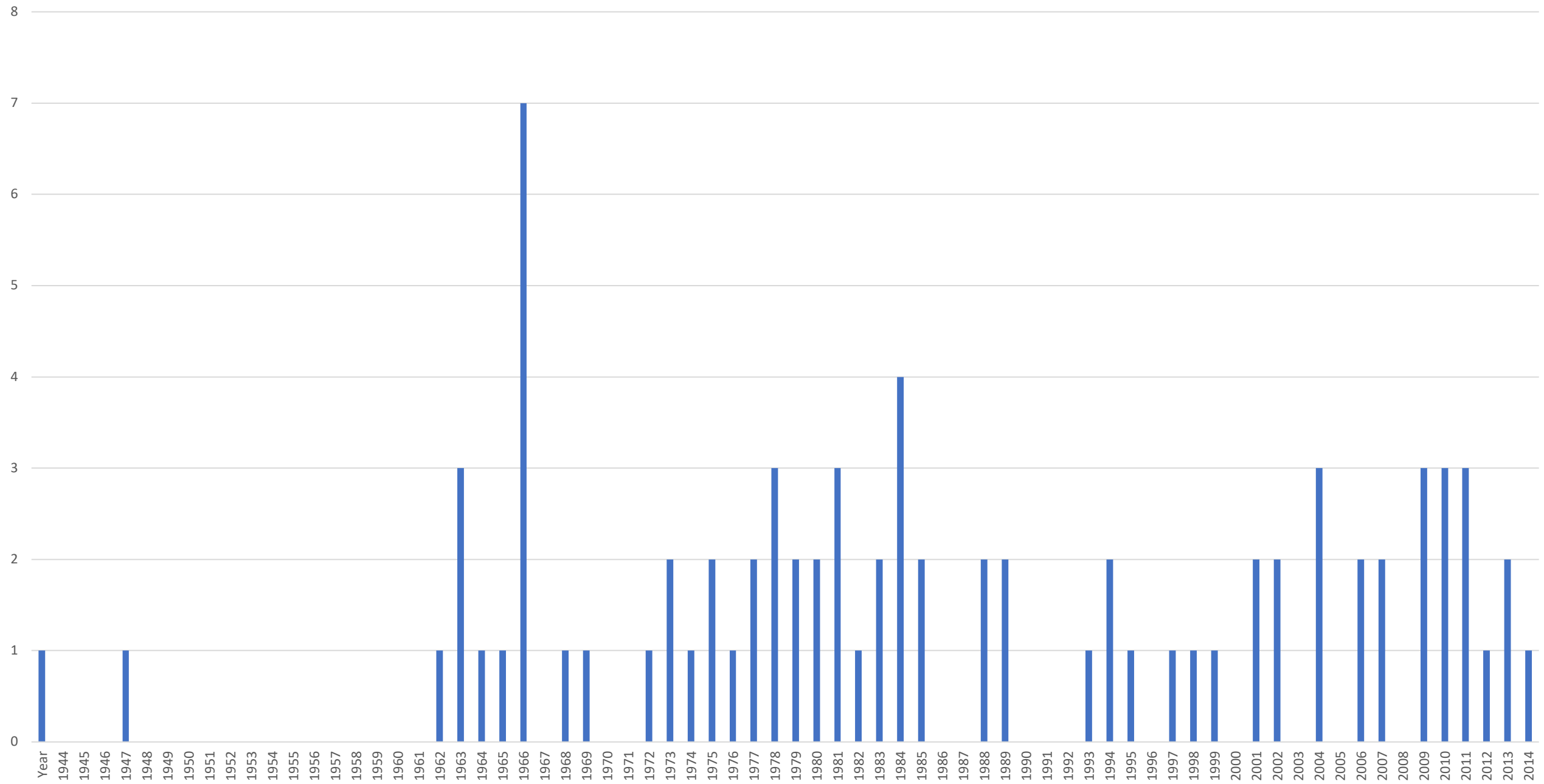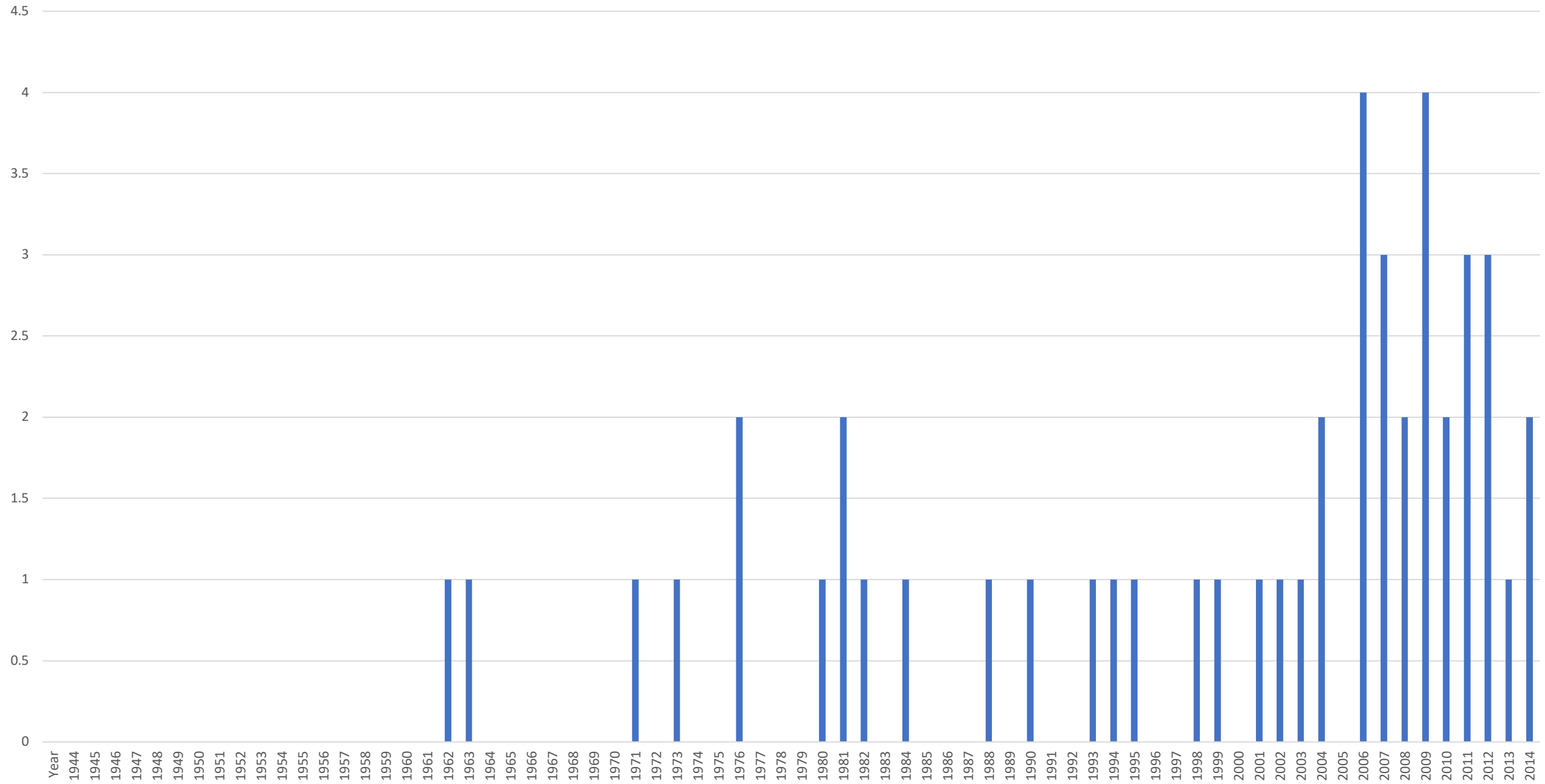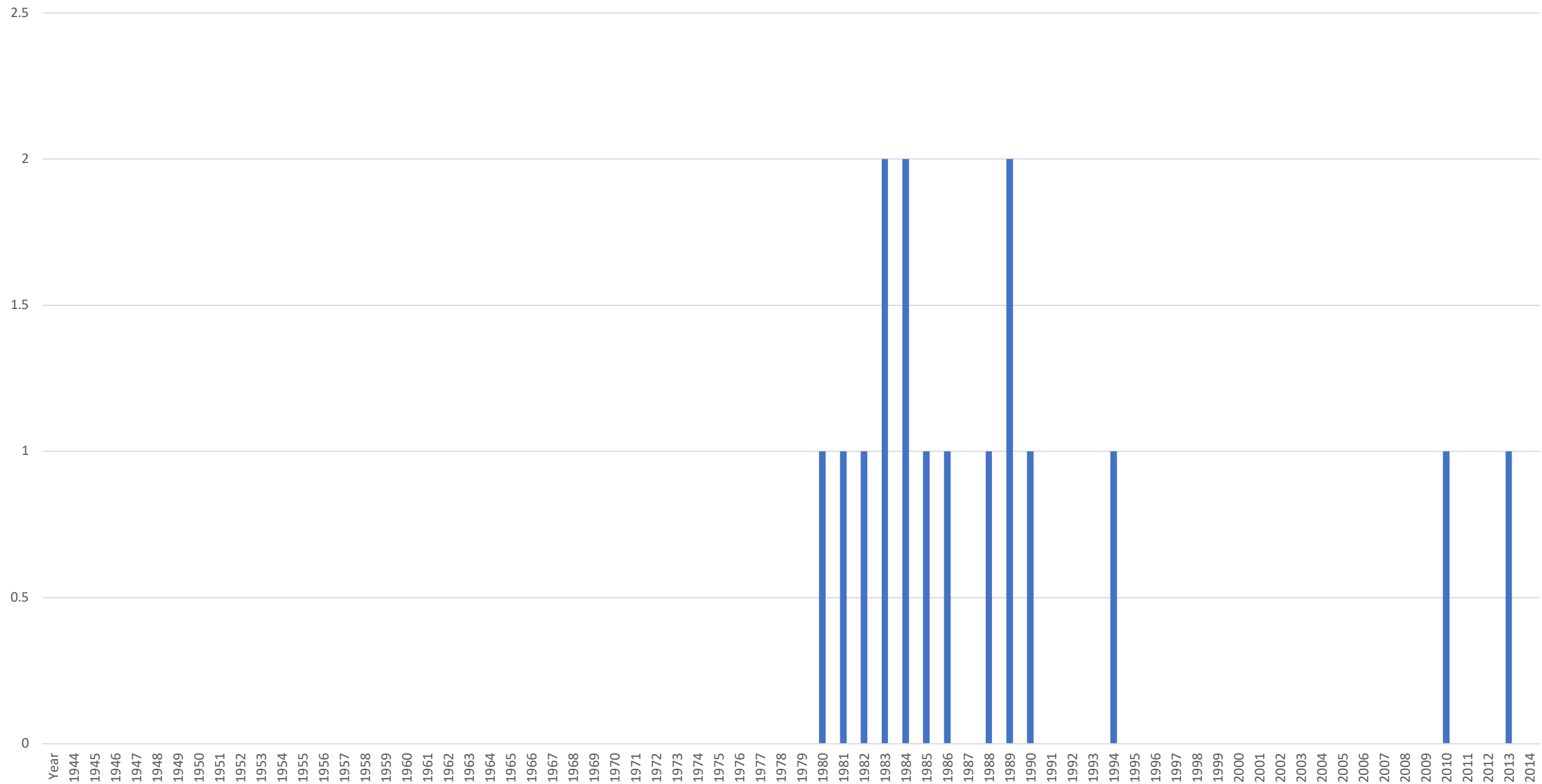
Text count

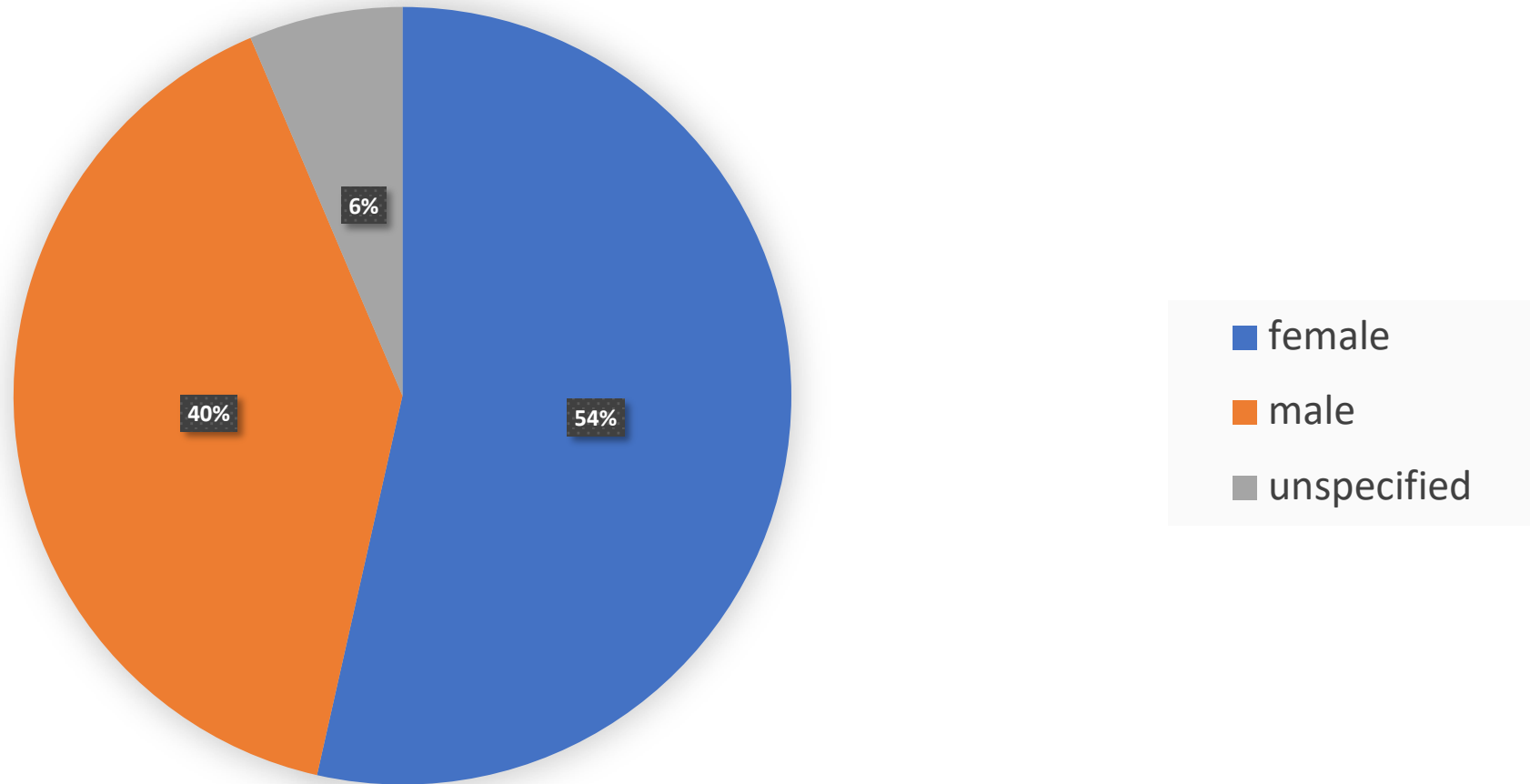Authentic Texts

# Syntax

Lexis

# Cohesion

# The Growth in Grammar Project

- Creating a corpus of educationally authentic writing from children in schools across England at Key Stages 1-4
- To be analysed for changes at the levels of lexis, phrase and clause.
  - NB: Analysis of structures used, not of accuracy.
- Corpus to be made publicly available from August 2018

# Corpus to date

| | Total texts | Total schools | Total children | Discipline | | |
|---|---|---|---|---|---|---|
| | | | | English | Hums | Science |
| Year 2 | 543 | 5 | 138 | 452 | 84 | 7 |
| Year 4 | 49 | 2 | 10 | 25 | 22 | 2 |
| Year 6 | 868 | 7 | 185 | 548 | 149 | 171 |
| Year 9 | 761 | 12 | 457 | 483 | 112 | 166 |
| Year 11 | 316 | 9 | 165 | 166 | 49 | 54 |
| All years | 2,537 | 23 | 955 | 1,674 | 416 | 400 |

# Corpus to date: Gender

# Corpus to date:
# English as an Additional Language

# Corpus to date:
# Socio-Economic Status

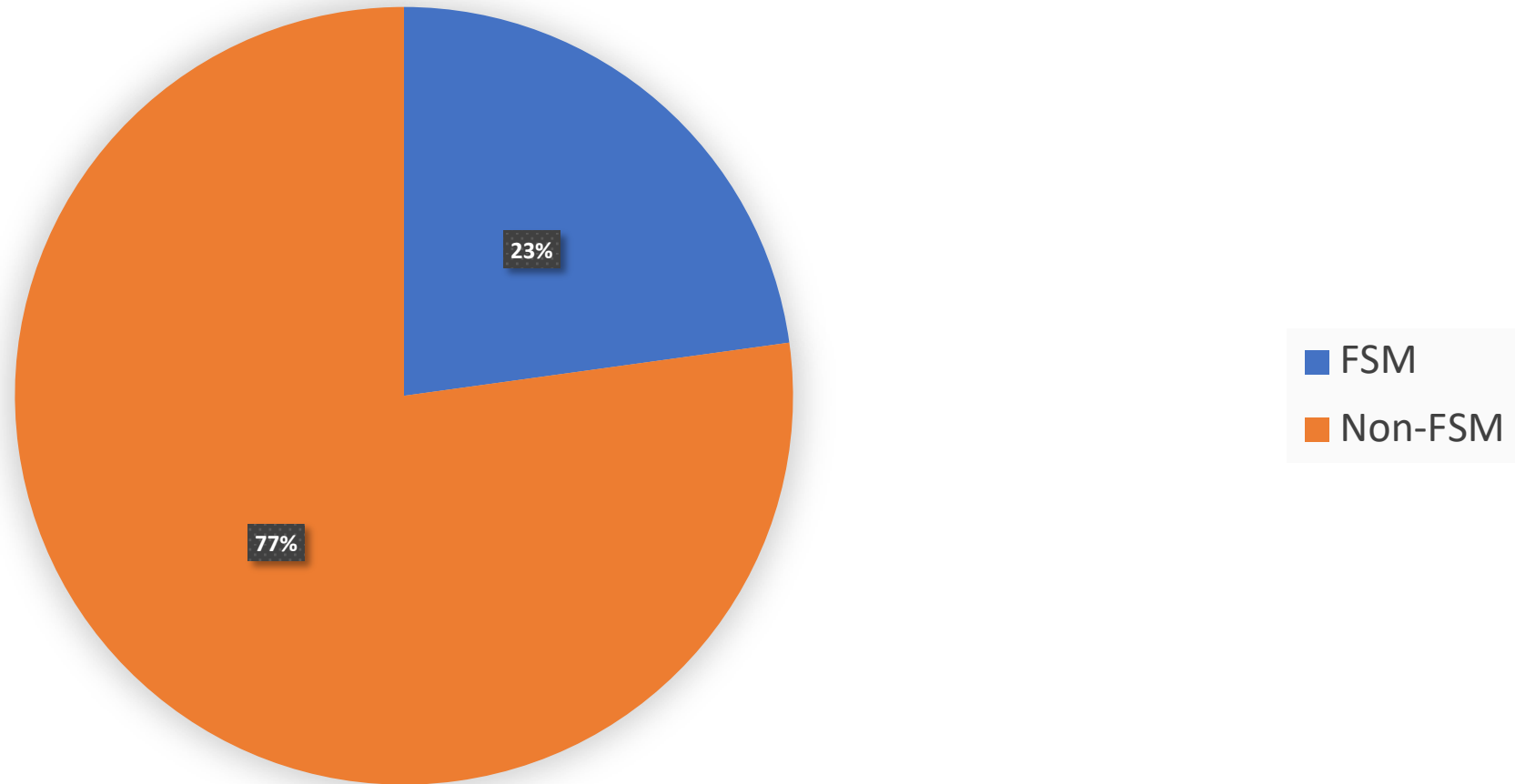# Ongoing:
# Coding Grammatical Features

```
                          Subordinate
                            Clauses

    Adverbial      Verb          Prepositional   Adjectival    Comparative
                   complement    complement      complement

       With          Subject
    subordinator

      Without         Object
    subordinator

      Fronted      Extraposed

      Relative      Adjective
                   complement
```

The captain ordered me to douse the fire in the galley in case the fire spread and burnt down the ship

The **captain** ordered **me** to douse the **fire** in the **galley** in case the **fire** spread and burnt down the **ship**

The captain ordered me to douse the fire in the galley in case the fire spread and burnt down the ship

det det det det det det

The **captain** ordered **me** to douse the **fire** in the **galley** in case the **fire** spread and burnt down the **ship**

The **captain** ordered **me** to douse the **fire** in the **galley** in case the **fire** spread and burnt down the **ship**

- the **captain**
- **me**
- the **fire** in the galley
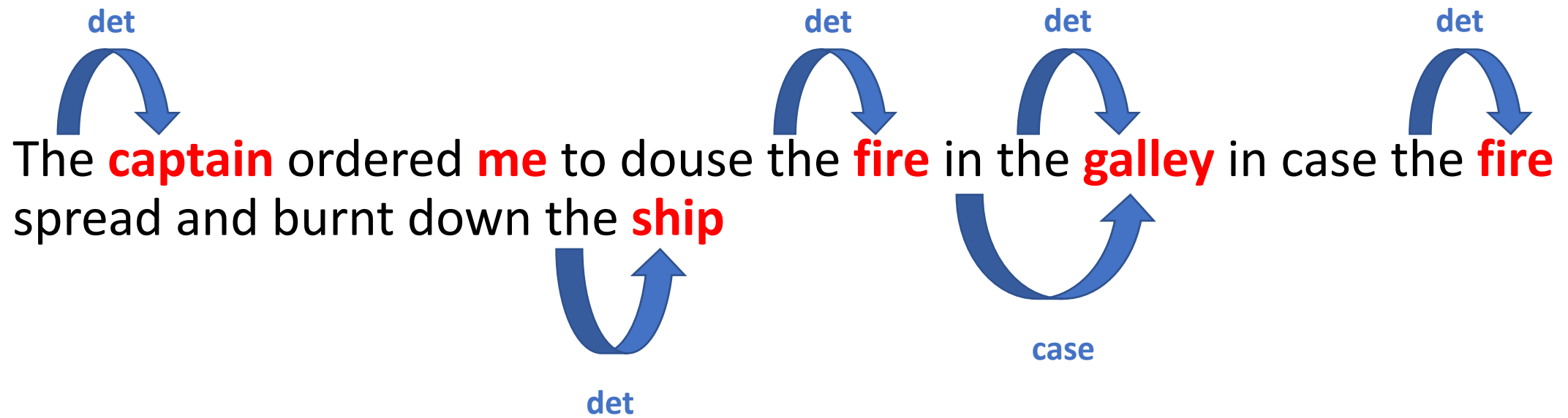  - in the **galley**
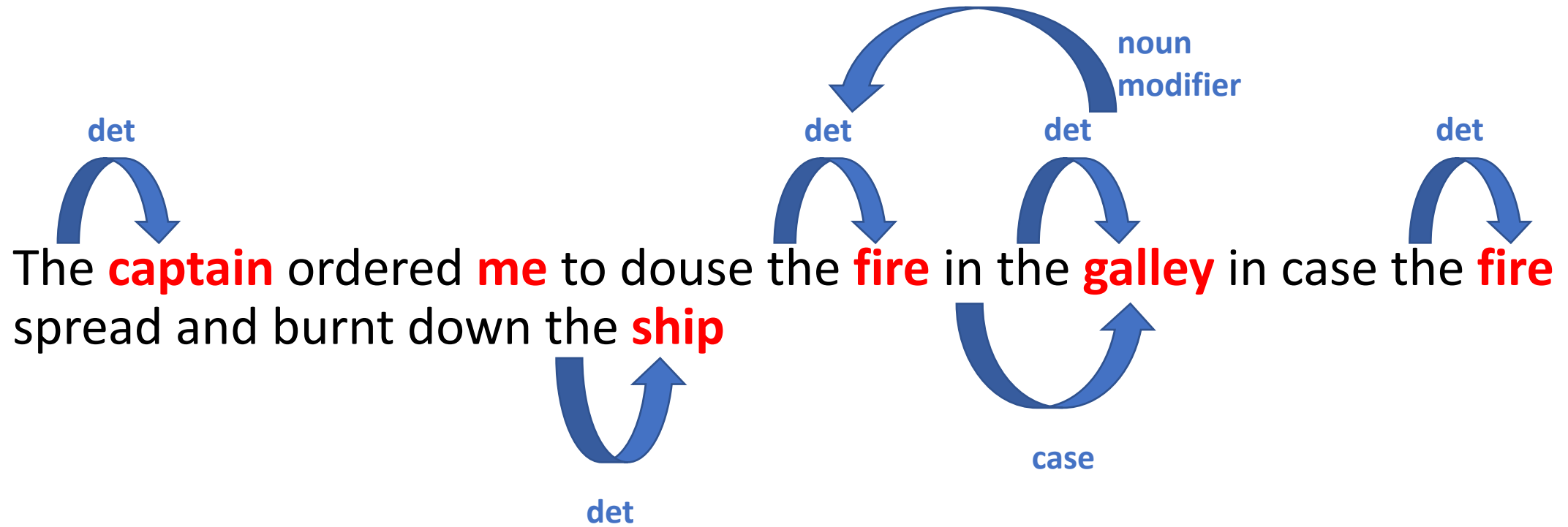- the **fire**
- the **ship**

The captain ordered me to douse the fire in the galley in case the fire spread and burnt down the ship

The **captain** ordered **me** to douse the **fire** in the **galley** in case the **fire** spread and burnt down the **ship**
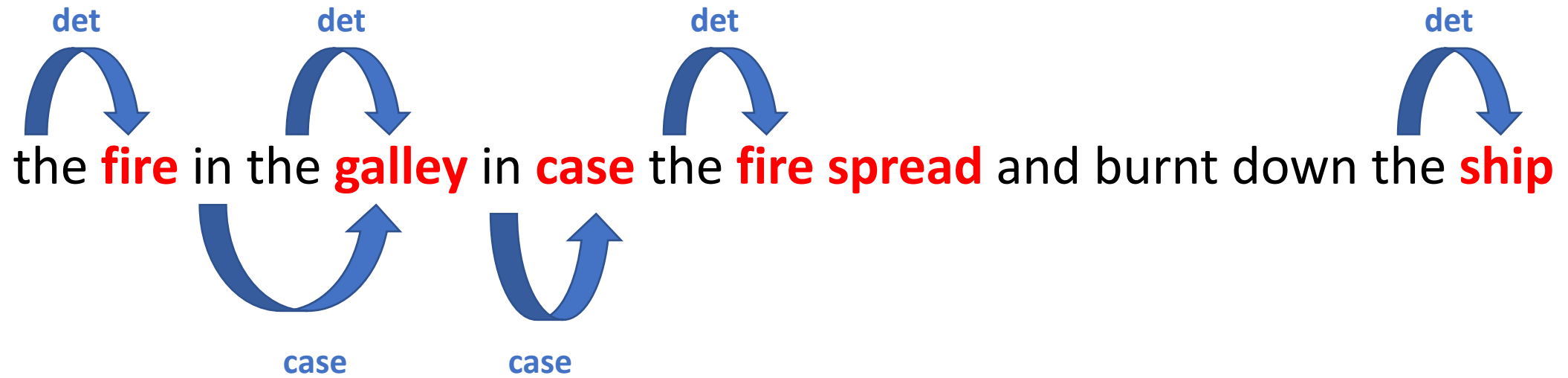
The **captain** ordered **me** to douse the **fire** in the **galley** in **case** the **fire** **spread** and burnt down the **ship**

the **fire** in the **galley** in **case** the **fire spread** and burnt down the **ship**

the **fire** in the **galley** in **case** the **fire spread** and burnt down the **ship**

the **fire** in the **galley** in **case** the **fire spread** and burnt down the **ship**
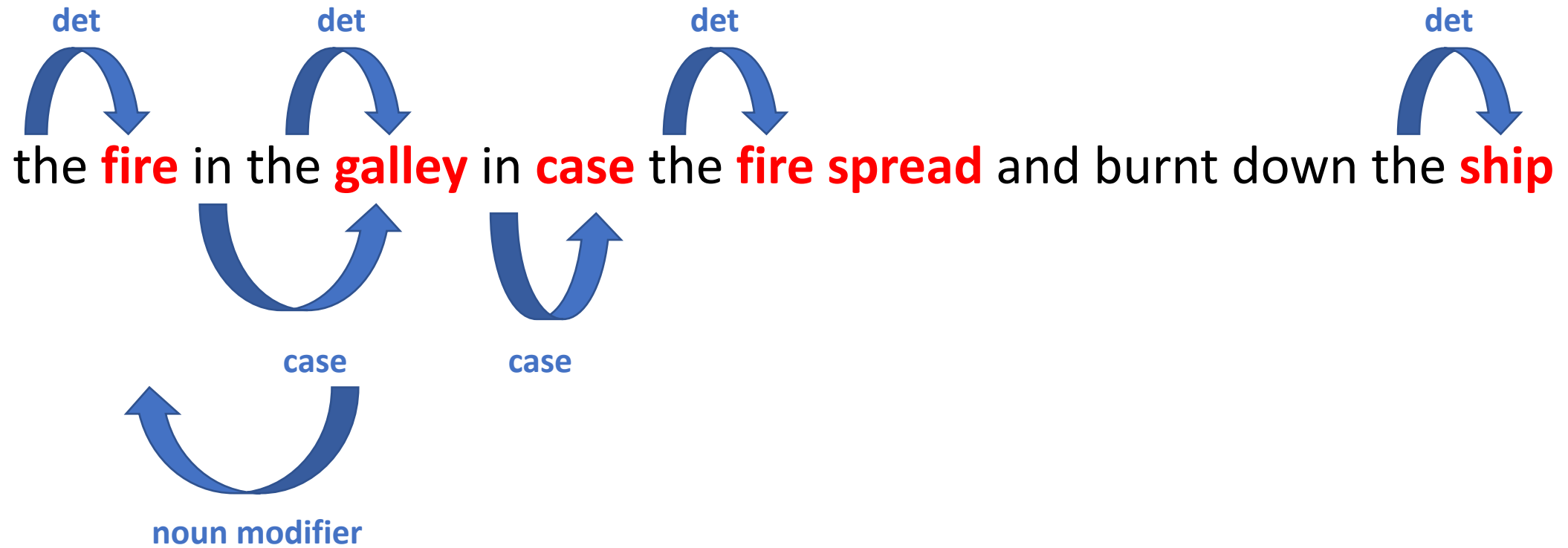
the **fire** in the **galley** in **case** the **fire spread** and burnt down the **ship**
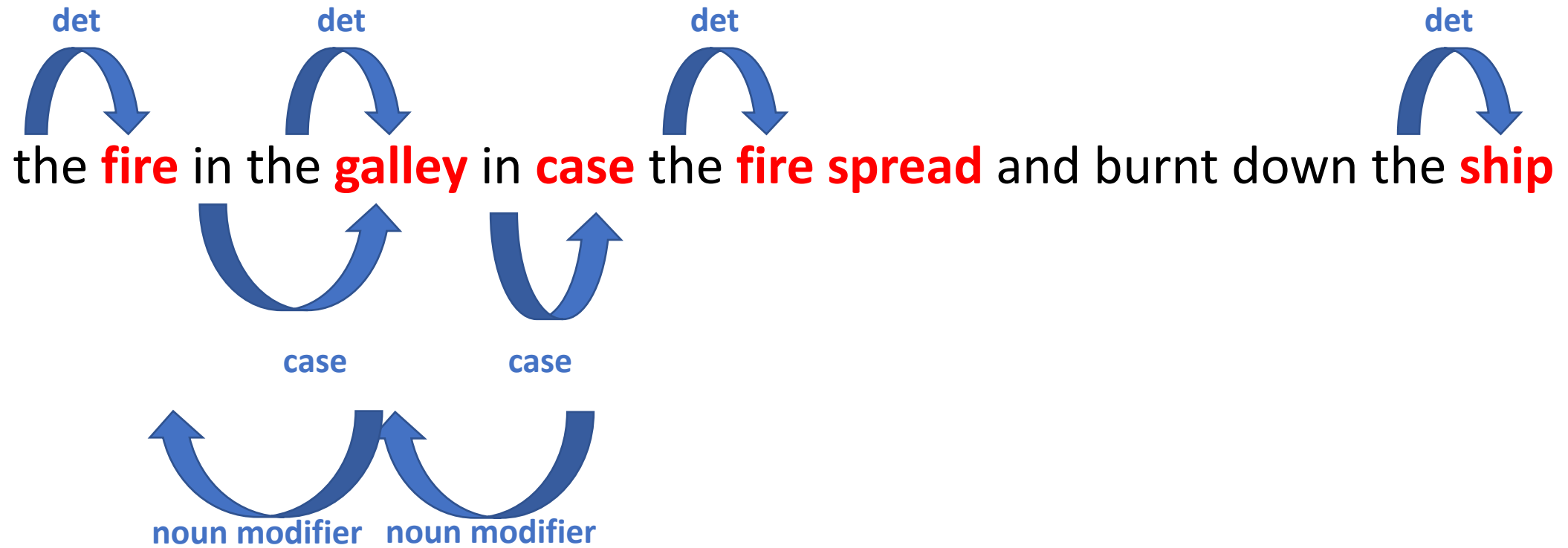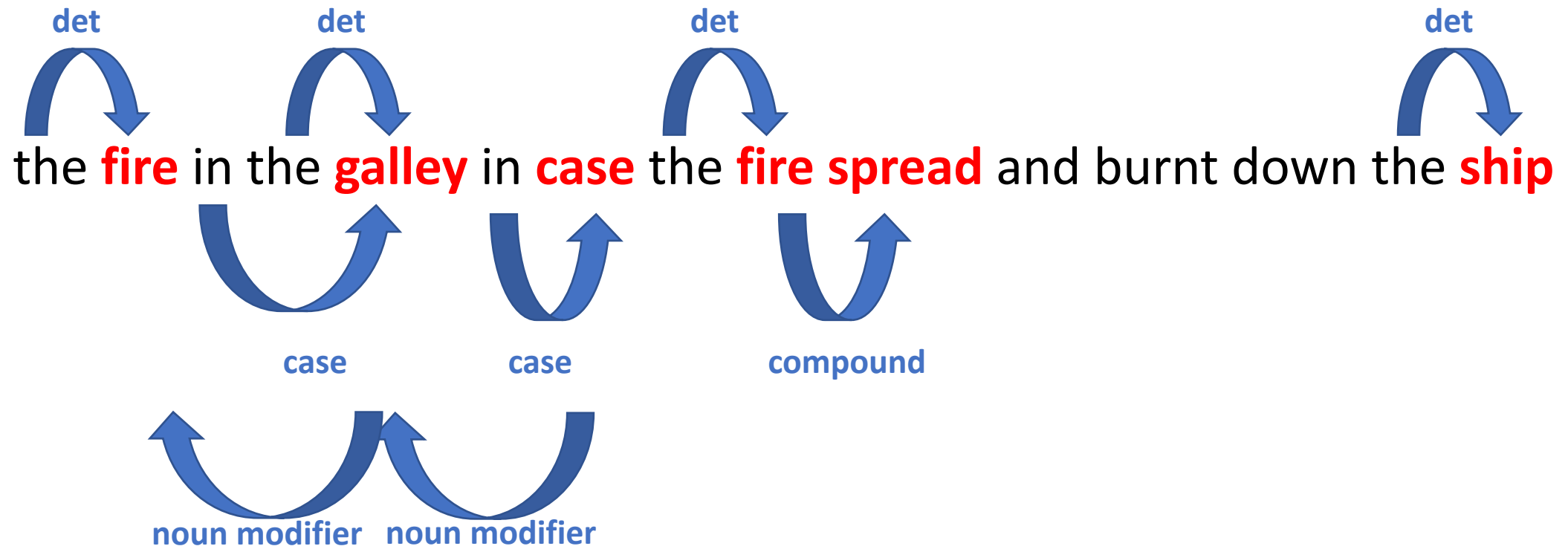
the **fire** in the **galley** in **case** the **fire spread** and burnt down the **ship**

- *the **captain***
- ***me***
- *the **fire** in the galley*
  - *in the **galley***
- *the **fire***
- *the **ship***

- *the **captain***
- **me**
- *the **fire** in the galley*
    - *in the **galley***
- *the **fire***
- *the **ship***

- *the **captain***
- **me**
- *the **fire** in the galley in case the fire spread*
    - *in the **galley** in case*
        - *in **case***
    - *the fire **spread***
        - *the **fire***
- *the **ship***

- 1 x 5-word NP
- 1 x 3-word NP
- 3 x 2-word NP
- 1 x 1-word NP

- Total NPs: 6
- Mean NP length: 2.5 words

- 1 x 10-word NP
- 1 x 5-word NP
- 1 x 3-word NP
- 4 x 2-word NP
- 1 x 1-word NP
- Total NPs: 8
- Mean NP length= 3.4 words

Mean Words per NP

Relative Clauses per Text

Relative Clauses per Text

*r* = .68

# Adverbial clauses per text



| | Yr 2 Lit | Yr 2 Non-lit | Yr 6 Lit | Yr 6 Non-lit | Yr 9 Lit | Yr 9 Non-lit | Yr 11 Lit | Yr 11 Non-lit |
|---|---|---|---|---|---|---|---|---|
| Manual | 3 | 0 | 12 | 3 | 10 | 5 | 20 | 5 |
| Stanford | 0 | 0 | 4 | 2 | 7 | 4 | 15 | 7 |

Adverbial clauses per text

*r* = .88

Mean Words per Adverbial Clause

# Mean Words per Adverbial Clause



$r = .82$

# Some Initial Findings: Vocabulary development

# Literature on word frequency and age

| Effect | Measure | Source | Ages |
|---|---|---|---|
| Decreases with age | % words < 1/100K | Finn 1977 | 9-10 < 16-17 |
| | % words not on high-frequency list | Olinghouse & Leaird 2009 | 7-8 > 9-10 |
| | % 1K words | Sun, Zhang & Scardamalia 2010 | 8 > 10 |
| | % 2K words | Sun, Zhang & Scardamalia 2010 | 8 < 10 |
| | % off-list words | Sun, Zhang & Scardamalia 2010 | 8 < 10 |
| Ambiguous age effects | % words not on high-frequency list | Lawton 1963 | 12 > 14 for working-class children, not middle class |
| | Mean frequency from reference corpus | Crossley et al 2011 | 14-15 < 18-19<br>14-15=16-17<br>16-17=18-19 |
| No age effects | P-Lex | Malvern et al 2004 | 7 = 11 = 14 |

# Sample for the current study

- Years 2, 6 and 9 only
- English/Humanities classes only
- Exclude texts with > 100 illegible words per 1,000
- Exclude poems
- Exclude samples more than 1SD from mean word length
- Spelling errors corrected
- Randomly select texts to give equal numbers in each year group

# Study Corpus

| | Schools | Writers | Texts | Text Length | | | | Genre | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | **Mean** | **Median** | **Min** | **Max** | **Story** | **Exposition** | **Persuasion** |
| Year 2 | 3 | 78 | 219 | 66.6 | 62 | 27 | 131 | 116 | 99 | 4 |
| Year 6 | 4 | 90 | 219 | 284.2 | 261 | 120 | 521 | 114 | 82 | 23 |
| Year 9 | 6 | 189 | 219 | 343.3 | 330 | 181 | 560 | 130 | 59 | 30 |

# Corpus structure

# TAALES* Indices

- Frequency
- Ngram frequency/association

# Frequency: 72 indices

- Range of reference corpora
- Separate indices for:
  - all words vs. content words vs. function words
  - raw frequency vs. log frequency

# Combining results from different corpora

| Sub-category (1) | Sub-category (2) | Cronbach's alpha | Deleted |
|---|---|---|---|
| All words | Raw | .99 | SUBTLEXus ($r$ = .41) |
| | Log | .98 | |
| Content words | Raw | .98 | |
| | Log | .98 | |
| Function words | Raw | .99 | SUBTLEXus ($r$ = .41) |
| | Log | .98 | |

|  |  | All words | | Content words | | Function words | |
|---|---|---|---|---|---|---|---|
|  |  | Raw | Log | Raw | Log | Raw | Log |
| All words | Raw | 1.00 |  |  |  |  |  |
|  | Log | 0.18 | 1.00 |  |  |  |  |
| Content words | Raw | -0.15 | 0.61 | 1.00 |  |  |  |
|  | Log | -0.27 | 0.80 | 0.76 | 1.00 |  |  |
| Function words | Raw | 0.83 | -0.17 | -0.30 | -0.42 | 1.00 |  |
|  | Log | 0.68 | -0.08 | -0.29 | -0.35 | 0.86 | 1.00 |

## Content Word Log Frequency

**Age:** Year 6: $b$=-.09, $t(78)$=-.49, $p$>.05; **Year 9: $b$=.56, $t(78)$=2.19, $p$<.05**
**Genre: non-literary: $b$=.56, $t(78)$=3.05, $p$<.05**
(random intercept for task)

## Function Word Log Frequency

**Age:** Year 6: $b$=-.15, $t(78)$=-.84, $p$>.05; **Year 9: $b$=-.72, $t(78)$=3.31, $p$<.05**
**Genre:** non-literary: $b$=.07, $t(78)$=0.43, $p$>.05
(random intercept for task)

# For example

One morning, five meerkats finished eating crunchy scorpians but there preditors like slithering snakes, lions and falcons. But one woke up and went meerkats, meerkats and woke up the others. One stepped on the snake and it went hiss hiss! And the other predators ran away. But a fennec fox was coming to take two pups. The scary fox wishes to eat them even the nice scorpians.

# Ngram measures: Proportion & Frequency

- Bigram & Trigram
- 7 x reference corpora; Association: 5 x reference corpora
- Proportion: 10K; 20K; 30K…100K)

# Combining results from different corpora

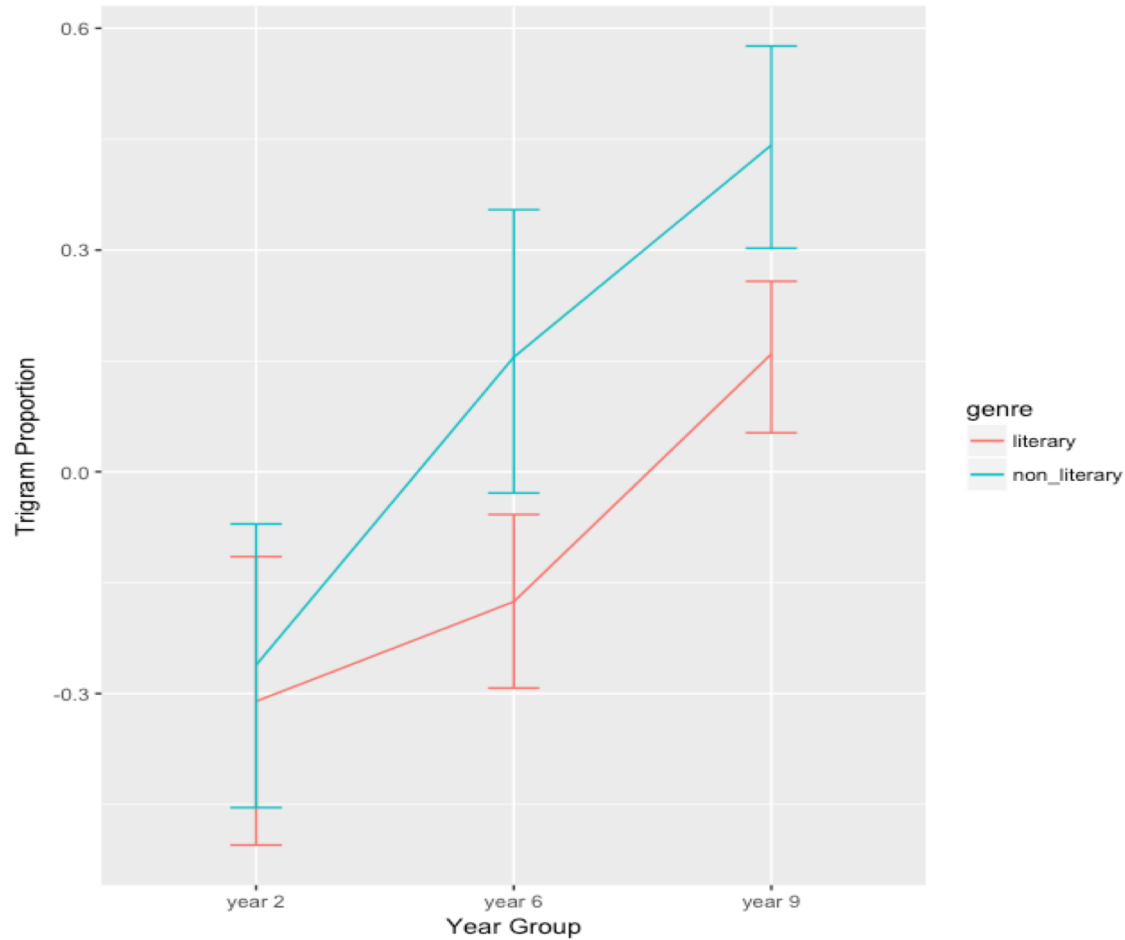| Category | Cronbach's alpha | Deleted |
|---|---:|---|
| Trigram Proportion (combines all frequency bands) | 1 | |
| Trigram Log Frequency | .93 | BNC Spoken: .14<br>BNC Written: .20<br>COCA Academic: .65 |

# Trigram Proportion

# Trigram Log Frequency

**Age:** Year 6: $b=.13$, $t(78)=.76$, $p>.05$; **Year 9: $b=.62$, $t(78)=2.06$, $p<.05$**
**Genre:** non-literary: $b=.19$, $t(78)=1.27$, $p>.05$
(random intercept for task)

**Age:** Year 6: $b=.13$, $t(78)=.1.05$ $p>.05$; Year 9: $b=.15$, $t(78)=-.07$, $p>.05$
**Genre:** non-literary: $b=.03$, $t(78)=.28$, $p>.05$
(random intercept for task)

# Trigram measures: Association

- Frequent n-grams tend to be combinations of frequent words:
  - *to be a*
  - *he didn't*
  - *back to the*
  - *it was the*
- Often use measures of *association*: is the combination frequent in relation to the frequency of the words?

# Trigram measures: Association

- *t-score:* Is the combination more frequent than chance would predict, given the frequency of the component words?

- *mutual information (MI):* how strongly do the component words predict each other?

- *delta-P:* conditional probability from one component to another

# Trigram measures: Association

- Range of reference corpora
- MI, MI2; t-score; Delta-P; Collexeme
- Trigram segmentation
  - Trigram 1 (*double* – *espresso please*)
  - Trigram 2 (*double espresso* – *please*)

# Combining results from different corpora

| Category | Cronbach's alpha | Deleted |
|---|---|---|
| MI | .94 | COCA Academic Tri 1: .58<br>COCA Academic Tri 2: .63 |
| MI2 | .97 | |
| t-score | .96 | COCA Academic Tri 1: .61<br>COCA Academic Tri 2: .66 |
| Delta-P | .91 | COCA Academic Tri 1: .63<br>COCA Academic Tri 2: .68 |
| Collexeme | .97 | COCA Academic Tri 1: .69<br>COCA Academic Tri 2: .63 |

| | MI | MI2 | t-score | Delta-P | Collexeme |
|---|---|---|---|---|---|
| MI | 1.00 | | | | |
| MI2 | 0.99 | 1.00 | | | |
| t-score | 0.66 | 0.66 | 1.00 | | |
| Delta-P | 0.42 | 0.44 | 0.42 | 1.00 | |
| Collexeme | 0.45 | 0.46 | 0.89 | 0.40 | 1.00 |

# Trigram associations

## Mutual Information

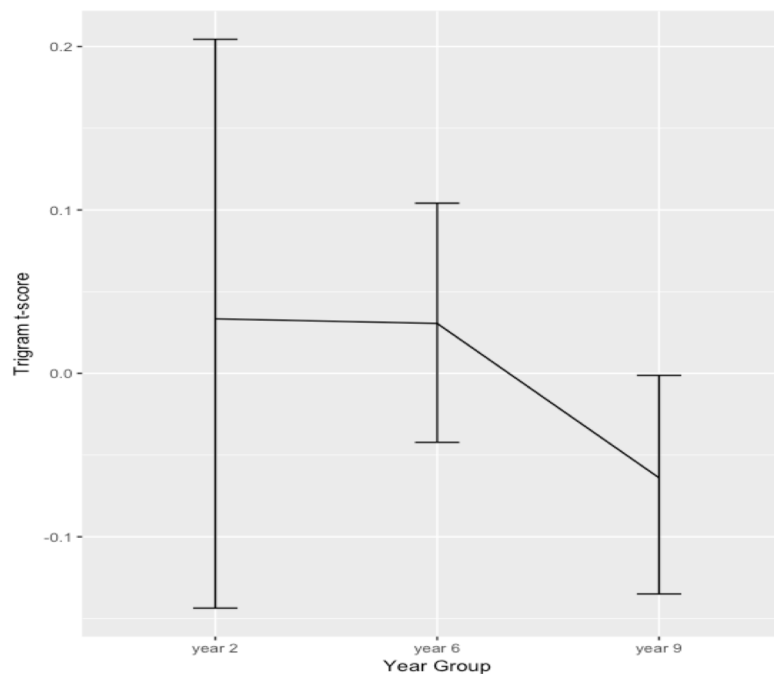

**Age: Year 6: *b*=.35, *t*(78)=3.01, *p*<.005;**
**Year 9: *b*=.38, *t*(78)=2.76, *p*<.01**
**Genre:** non-literary: *b*=.02, *t*(78)=.17, *p*>.05
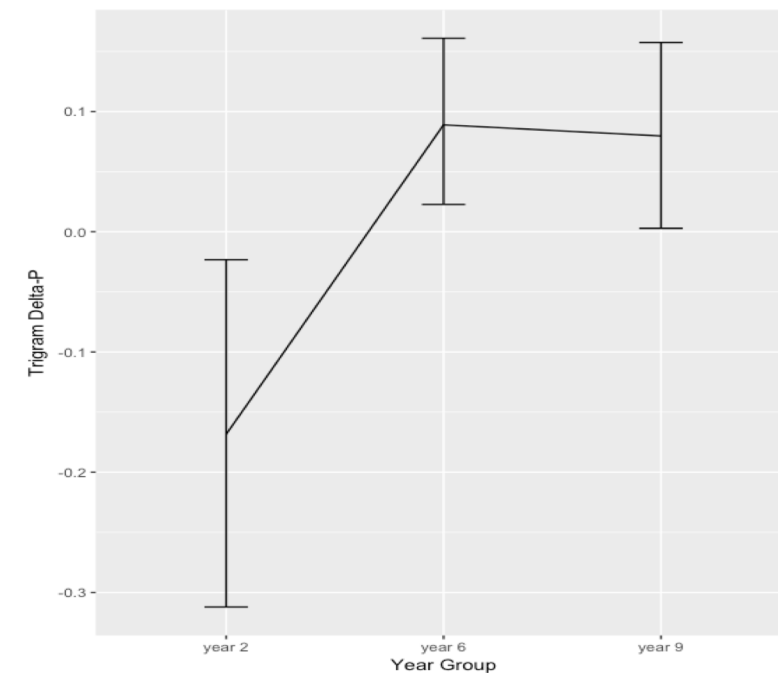(random intercept for task)

## t-score



**Age:** Year 6: *b*=.03, *t*(78)=.28, *p*>.05;
Year 9: *b*=-.1, *t*(78)=-.80, *p*>.05
**Genre:** non-literary: *b*=.00, *t*(78)=-.03, *p*>.05
(random intercept for task)

## Delta-p



**Age: Year 6: *b*=.40, *t*(78)=3.35, *p*<.005;**
**Year 9: *b*=-.32, *t*(78)=2.35, *p*<.05**
**Genre:** non-literary: *b*=-.14, *t*(78)=-1§.39, *p*>.05
(random intercept for task)

# Conclusions - methodological

- Counts from different reference corpora mostly consistent
- Log frequencies enable patterns to emerge more clearly

# Conclusions: word frequency

- Mean content word frequency increases with age
- Mean function word frequency decreases with age

# Conclusions: trigrams

- Percentage of ngrams attested in corpora increases with age
- MI & DP of attested trigrams increase with age

So...

- Corpus research under-exploited in study of L1 English writing development.
- Our corpus to be completed early 2018
- Online late 2018
- Full analyses of vocabulary, NP-expansion, subordination soon.
- And the book of the lit. review…
- Future prospects:
  - Historical oral/written corpus
  - Studies of attainment

# Keep in touch!

- Twitter:      @growing grammar
- Facebook:  facebook.com/growthingrammar
- Email:          p.l.durrant@exeter.ac.uk

# References

- Crossley, S. A., Weston, J. L., Sullivan, S. T. M., & McNamara, D. (2011). The development of writing proficiency as a function of grade level: A linguistic analysis. *Written Communication, 28*, 282-311.

- Department of Education (2014). *The national curriculum in England: Framework document. December 2014.*

- Finn, P. J. (1977). Computer-aided description of mature word choices in writing. In C. R. Cooper & L. Odell (Eds.), *Evaluating writing: Describing, Measuring, Judging* (pp. 69-89). Urbana, Illinois: National Council of Teachers of English.

- Malvern, D., Richards, B. J., Chipere, N., & Duran, P. (2004). *Lexical diversity and language development*. Basingstoke: Palgrave Macmillan.

- Olinghouse, N., G., & Leaird, J. T. (2009). The relationship between measures of vocabulary and narrarive writing quality in second- and fourth-grade students. *Reading and Writing, 22*, 545-565.

- Sun, Y., Zhang, J., & Scardamalia, M. (2010). Knowledge building and vocabulary growth over two years, Grades 3 and 4. *Instructional Science, 38*, 147-171.