

Slipping through the cracks of e-lexicography: lessons from

COLLOcaid

Ana Frankenberg-Garcia

Geraint Rees

Robert Lew

Jonathan C. Roberts

Nirwan Sharma

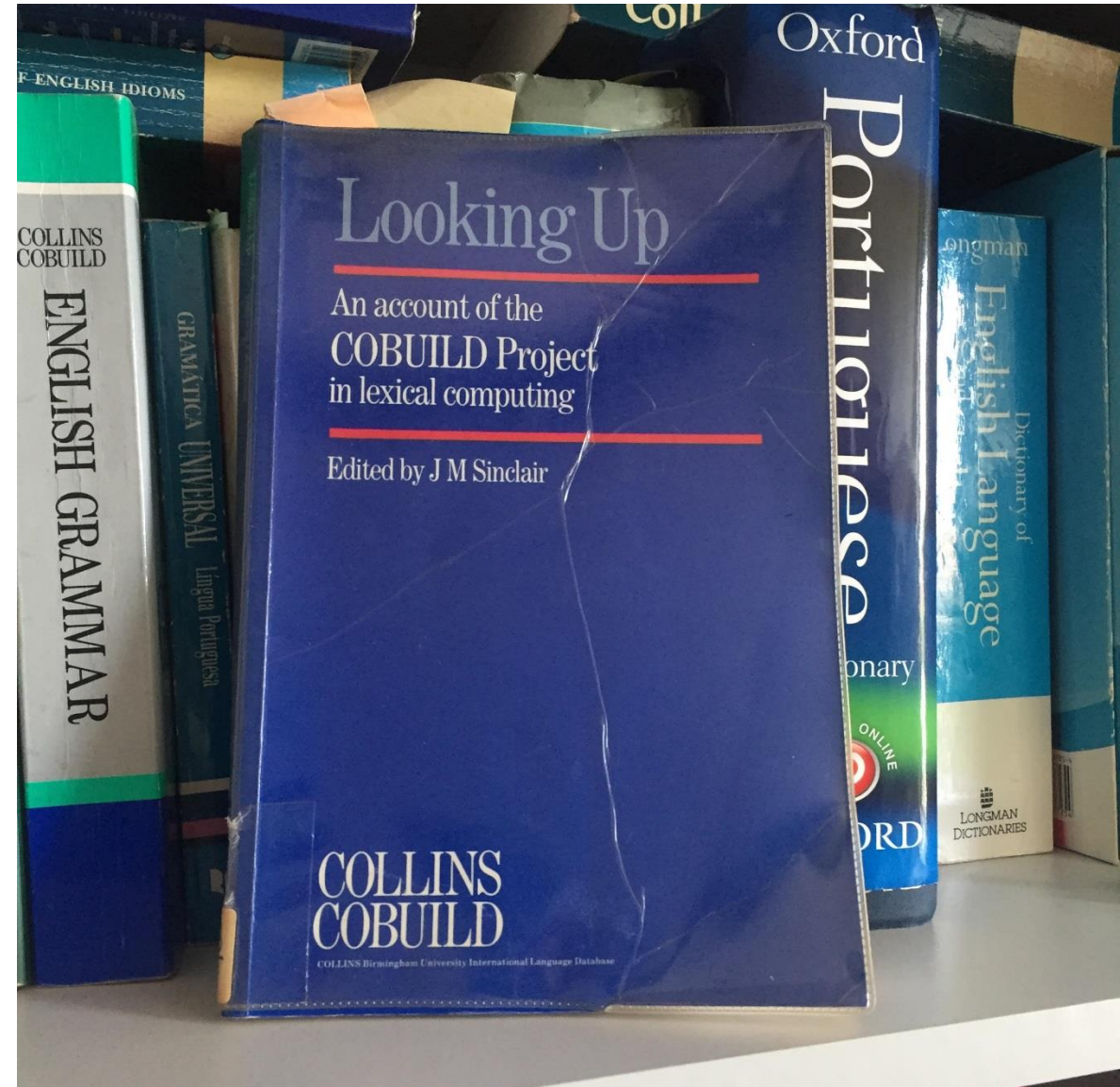
Peter Butcher



1987

First developments in e-lexicography

- COBUILD project (Sinclair 1987)
- Lexical analyses aided by corpora
 - Concordances
 - Word lists
 - Collocations



Things have evolved...

Word Sketches

- Lexical summary of how words are used with other words (Kilgarriff et al. 2004, 2014)

logDice

- Better statistics for corpus-based lexicography (Rychlý 2008, Gablasova et al. 2017, Frankenberg-Garcia 2018)

GDEX

- Easier to find good dictionary examples from corpora (Kilgarriff et al. 2008)

e-Lexicography today

- Headword lists & defining vocabularies adjusted to better reflect language use
- Definitions enhanced with linguistic data beyond introspection
- Improved coverage of syntactic patterns, lexical collocations & phraseology
- Abundancy of authentic examples of words in context

And yet...



- When compiling the lexical database for ColloCaid, we noted...
- Problems that seem to be slipping through the cracks of e-lexicography

We present next lessons learnt and some examples

What is ColloCaid?

Academic writing is hard...

collocaid



ColloCaid is a text editor that assists writers with academic English collocations

How does ColloCaid work?

The screenshot illustrates the ColloCaid interface, which is designed to assist with word choice in a word processing application. It features a standard menu bar (File, Edit, Insert, View, Format, Table, Tools, Help) and a paragraph toolbar with various formatting options. The main text area shows the word "advantage" highlighted in blue. A context menu is open over the word, displaying a list of suggestions:

- take advantage V+ ▶
- selective advantage Adj+ ▶
- advantage of (method/approach) ▶
- advantage over (instruments/time/rivals) ▶
- selective advantage ▶
- main advantage ▶
- mutual advantage ▶
- added advantage ▶
- obvious advantage ▶
- potential advantage ▶
- unfair advantage ▶
- distinct advantage ▶
- More

The "More" option at the bottom of the context menu is circled in red. To the right of the main text area, a vertical list of suggestions is visible, including:

- selective advantage ▶
- mutual advantage ▶
- unfair advantage ▶
- distinct advantage ▶
- main advantage ▶
- obvious advantage ▶
- added advantage ▶
- potential advantage ▶
- clear advantage ▶
- great advantage ▶
- strategic advantage ▶
- relative advantage ▶
- major advantage ▶
- significant advantage ▶
- absolute advantage ▶
- full advantage ▶
- considerable advantage ▶
- additional advantage ▶
- unique advantage ▶

The interface also shows a second instance of the word "advantage" in a different paragraph, and a third instance where the word "Another" is followed by "advantage".

What's different about ColloCaid?

- Can help writers expand their academic collocation repertoire
 - Real-time collocation reminders
 - Writers can access collocations they might otherwise not remember to use
- Interactive text editor integration
 - Writers don't have to stop writing
 - Collocations suggestions are easy to find
 - But only given when sought, so as not to disrupt writing
- Data-driven learning (Johns 1991)
 - Suggestions shown, not explained
 - Metalanguage kept to a minimum
- Curated lexical data, so that users don't get distracted with
 - Irrelevant or misleading information
 - An overload of information

What lessons did we learn?

Lessons from ColloCaid



1. Issues with node selection
2. Issues with collocation research
3. Issues with corpus examples

Node selection

Node selection

() *Where in my text can I fit in **conduct research**?*

(x) *What verb can I use with **research**?*

- Starting point is the node, not the collocation
- 500 maximally useful collocation nodes for general academic English
 - 36-month project
- Nouns, verbs and adjectives from 3 well-established word lists of general academic English
 - Extracted using different methods and different corpora
- Build on the strengths of each list

Node selection

1. **AKL** - Academic Keyword List (Paquot 2010)

- Emphasis on academic keyness

2. **AVL-BAWE** - Subset of Academic Vocabulary List (Gardner & Davies 2014) that overlaps with BAWE word list (Durrant 2016)

- Emphasis on novice academic writing vocabulary

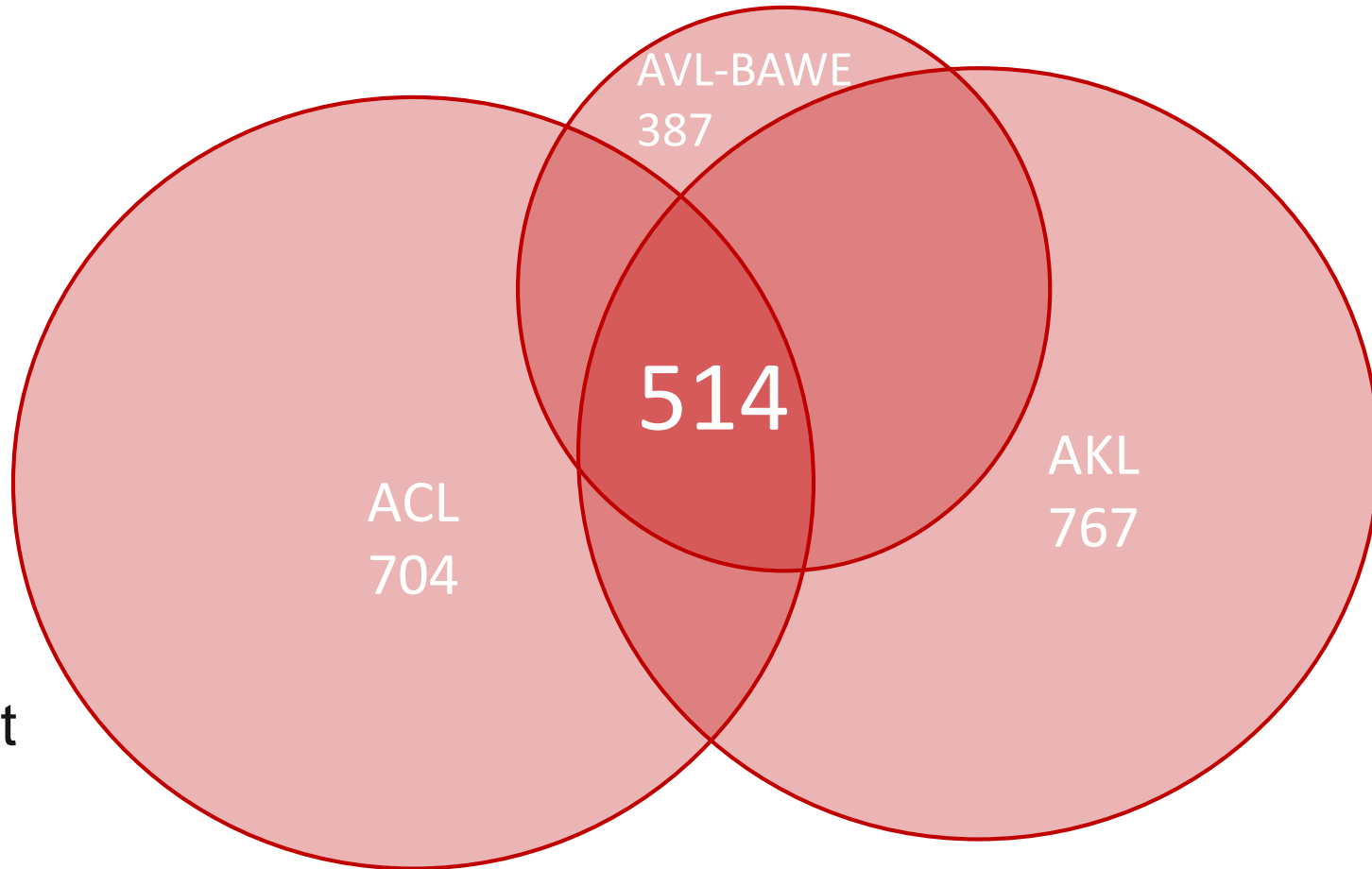
3. **ACL** - Academic Collocations List (Ackermann & Chen 2013), using headwords in appendix to Longman Collocations Dictionary (Mayor 2013)

- e.g. for *conduct research*, headword = *research*
- Emphasis on strong collocations

Node selection

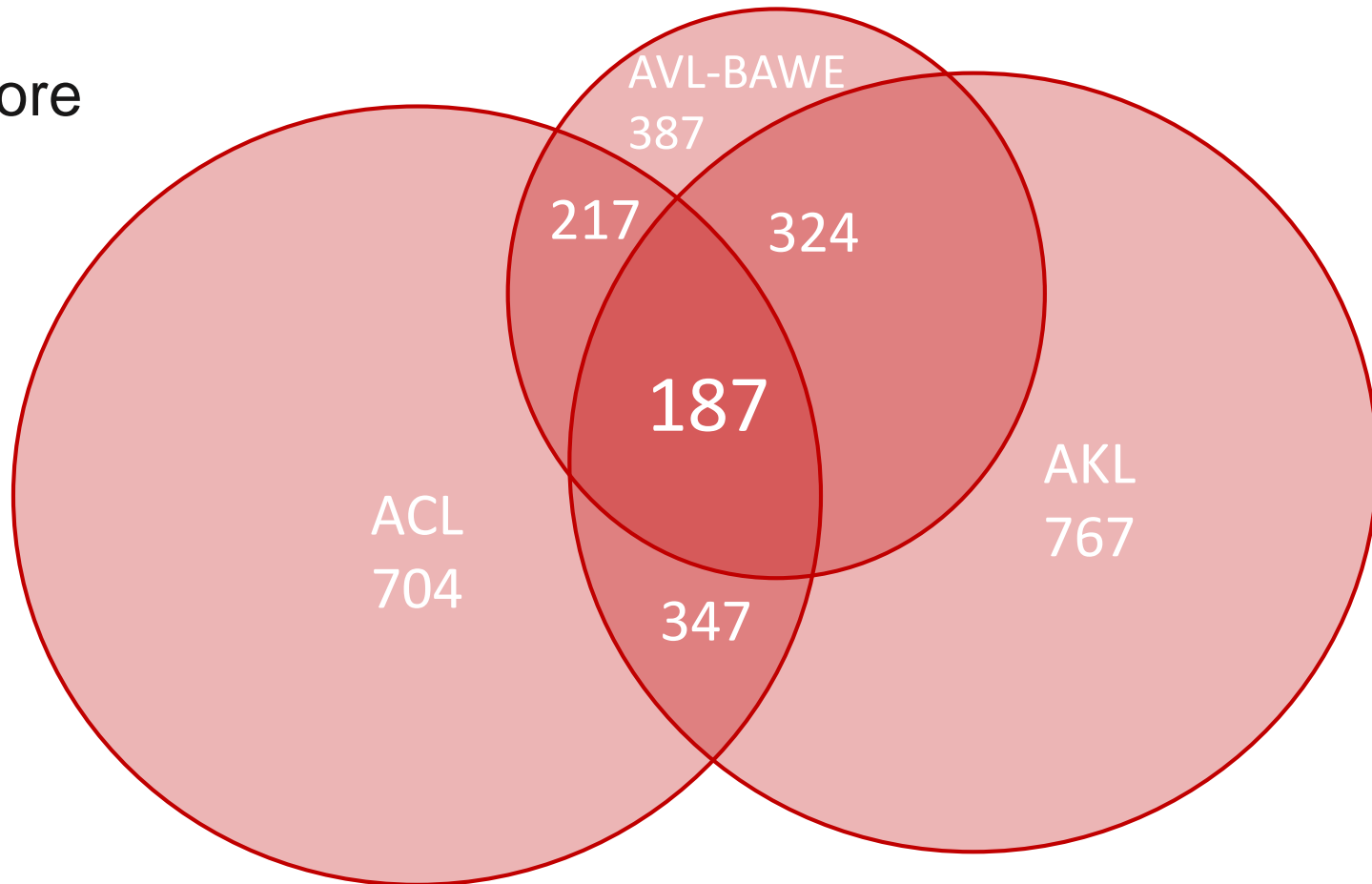
Word list intersections

- AKL (biggest, academic keyness)
 - AVL-BAWE (smallest, students)
 - ACL (biggish, strong collocations)
-
- 514 lemmas overlapping in at least two lists
 - Perfect fit for project scope 😊



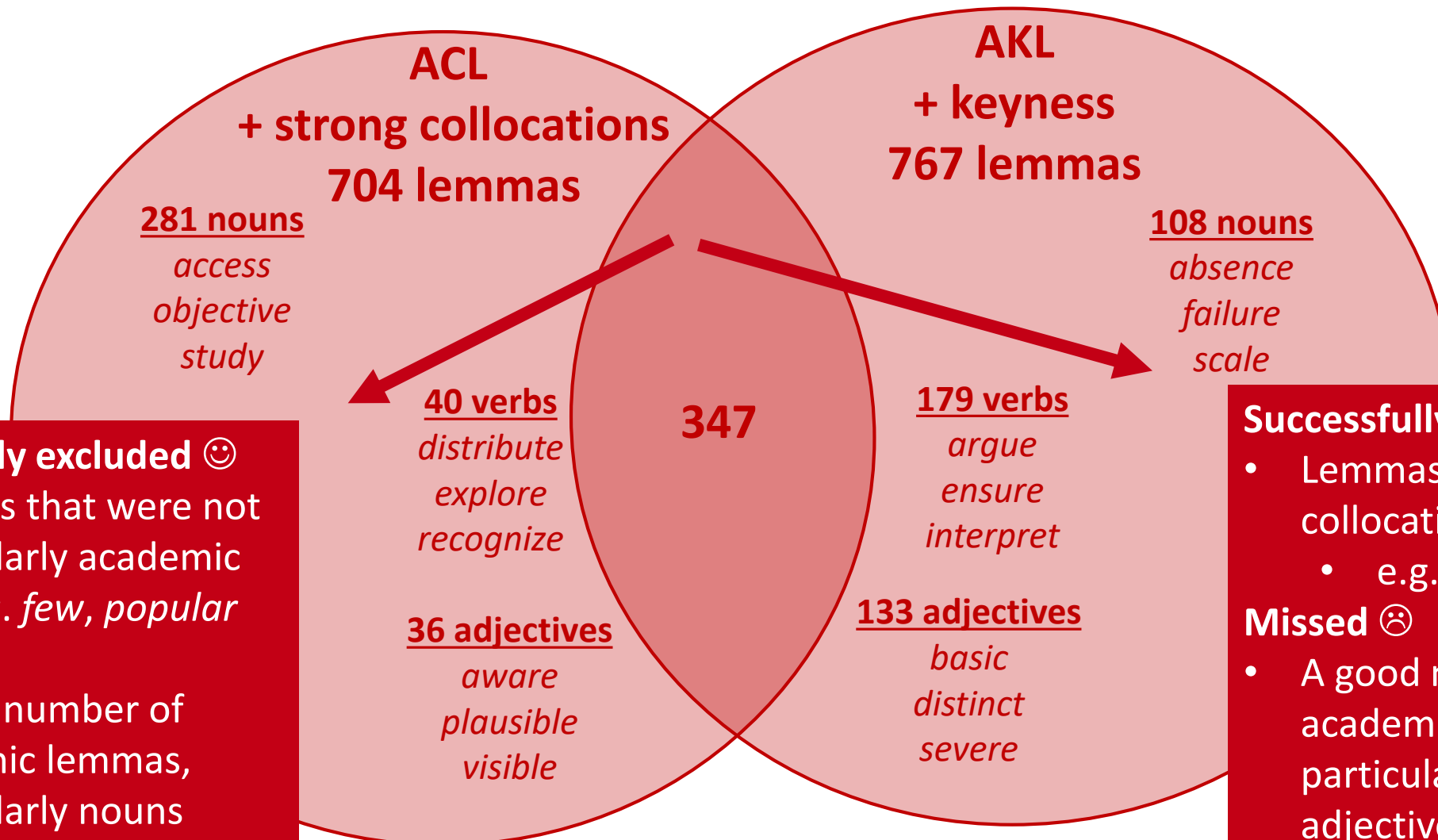
Node selection issues

- Even though all 3 lists are about core general academic English vocabulary...
- They overlap less than expected!



Words lists aiming for similar vocabulary coverage (i.e. core academic vocabulary) can differ substantially depending on corpora and extraction methods used

Node selection issues



Successfully excluded 😊

- Lemmas that were not particularly academic
 - e.g. *few*, *popular*

Missed 😞

- A good number of academic lemmas, particularly nouns

Successfully excluded 😊

- Lemmas that were not collocationally productive
 - e.g. *actual*, *prime*

Missed 😞

- A good number of academic lemmas, particularly verbs & adjectives

Node selection issues

AVL-BAWE
+ novice academic
387 lemmas

33

Successfully excluded 😊

A few less relevant lemmas:

- *university* > Cambridge University Press, etc.
- non-gradable adjectives: *existing*, *economic*, *current*, etc.

Missed 😞

Lexical teddy bears? (Hasselgren 1994)

Genre?

But surely some are relevant, e.g.:

- *timely/systematic/satisfactory manner*;
- *shown/summarized/presented in Table x*;
- *slightly/substantially/suitably modified*;
- *mutually/highly/especially beneficial*.

9 nouns

appendix
century
flow
inclusion
manner
project
region
table
university

9 verbs

address
calculate
estimate
justify
modify
observe
range
situate
test

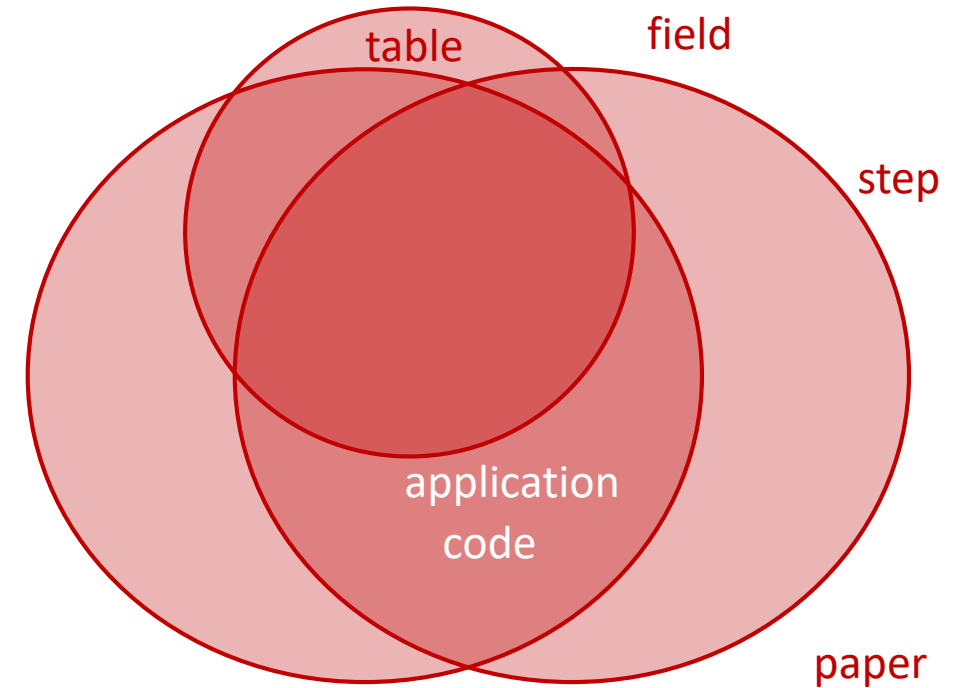
15 adjectives

above
accurate
beneficial
broad
continuous
current
detailed
economic
existing
external
given
global
increased
numerous
varying

Node selection issues

We also noted

1. Basic lemmas with major academic meanings were missing from well-known academic word lists
2. Less central lemmas appeared to be overrated in comparison



Node selection issues

Why were basic academic lemmas with major academic meanings missed?

- Words lists used disregard sense distinctions
- Less sensitive to academic senses of lemmas widely used in general language



Important academic senses get diluted and escape frequency thresholds

- Ideally, most obvious missing academic lemmas need to be added

Node selection issues

Why did less central academic lemmas seem overrated in comparison?

Two telling examples:

1. application

= **request** (*a successful patent application*)

= **software** (*IT applications require significant innovative effort*)

= **use** (*the application of statistical techniques to test assumptions*)

Lemmas with distinct general academic senses lumped together and jump rank queue

2. code (not interdisciplinary, but multidisciplinary)

In Computing

In Biology

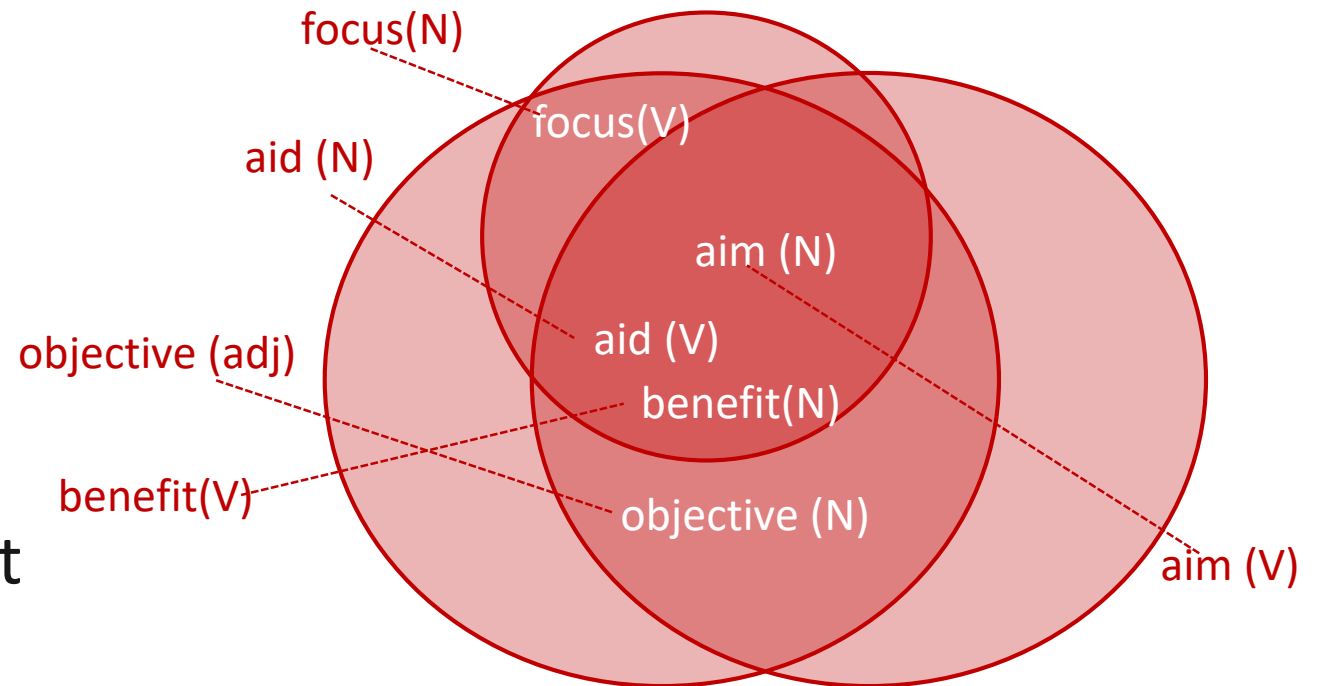
In Language and Linguistics

Terms with distinct discipline-specific senses lumped together and also jump rank queue

Node selection issues

Another issue was word class

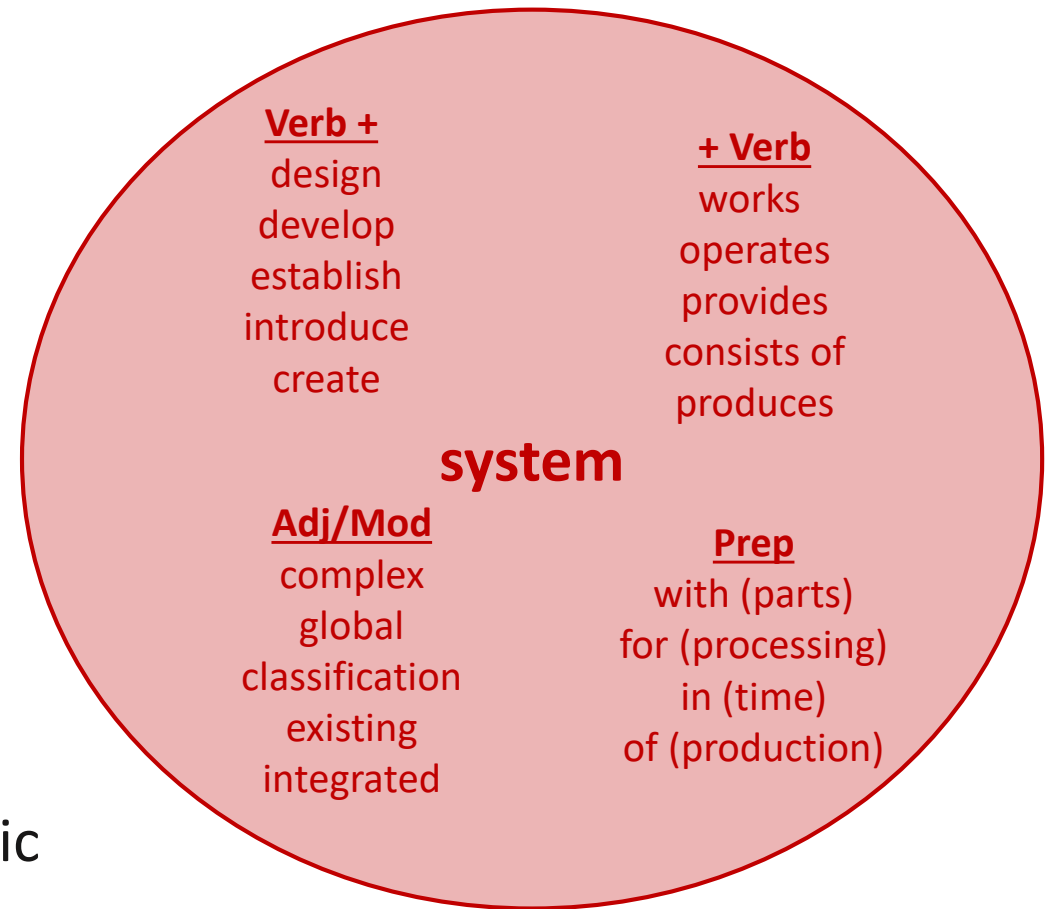
- Some homographs are naturally more frequent in one class than another
- But this can give the impression that lexical coverage is inconsistent
- Decision to include homographs
 - Even if they did not meet our thresholds



Collocation research

Collocation research

- Collocations likely to be looked up for each node
- Interdisciplinary lexical collocations
- Grammatical collocations (syntactic patterns)
- Expert academic English corpora
 - **Oxford Corpus of Academic English**
 - Pearson International Corpus of Academic English
 - COCA academic



Collocation research

Word sketches (Kilgarriff et al. 2004, 2014)

- Collocates sorted per grammar relation

research + ADJ

- logDice & co-occurrence thresholds
- excluded discipline-specific words (dispersion)
- excluded open-choice combinations
 - e.g. *more research*

qualitative	1,522	10.51	...
qualitative research			
future	1,203	9.89	...
for future research			
quantitative	751	9.53	...
quantitative research			
further	873	9.22	...
further research			
previous	693	9.04	...
previous research			
empirical	523	8.91	...
empirical research			
recent	544	8.52	...
recent research			
market	472	8.38	...
market research			

Collocation research issues

- Is anything missing?
- *carry out + research* (+5, -5) = 303!
- Space between *carry* and *out* is problematic
- Further examples
 - *set out*
 - *set up*
 - *at stake*
 - *at hand*

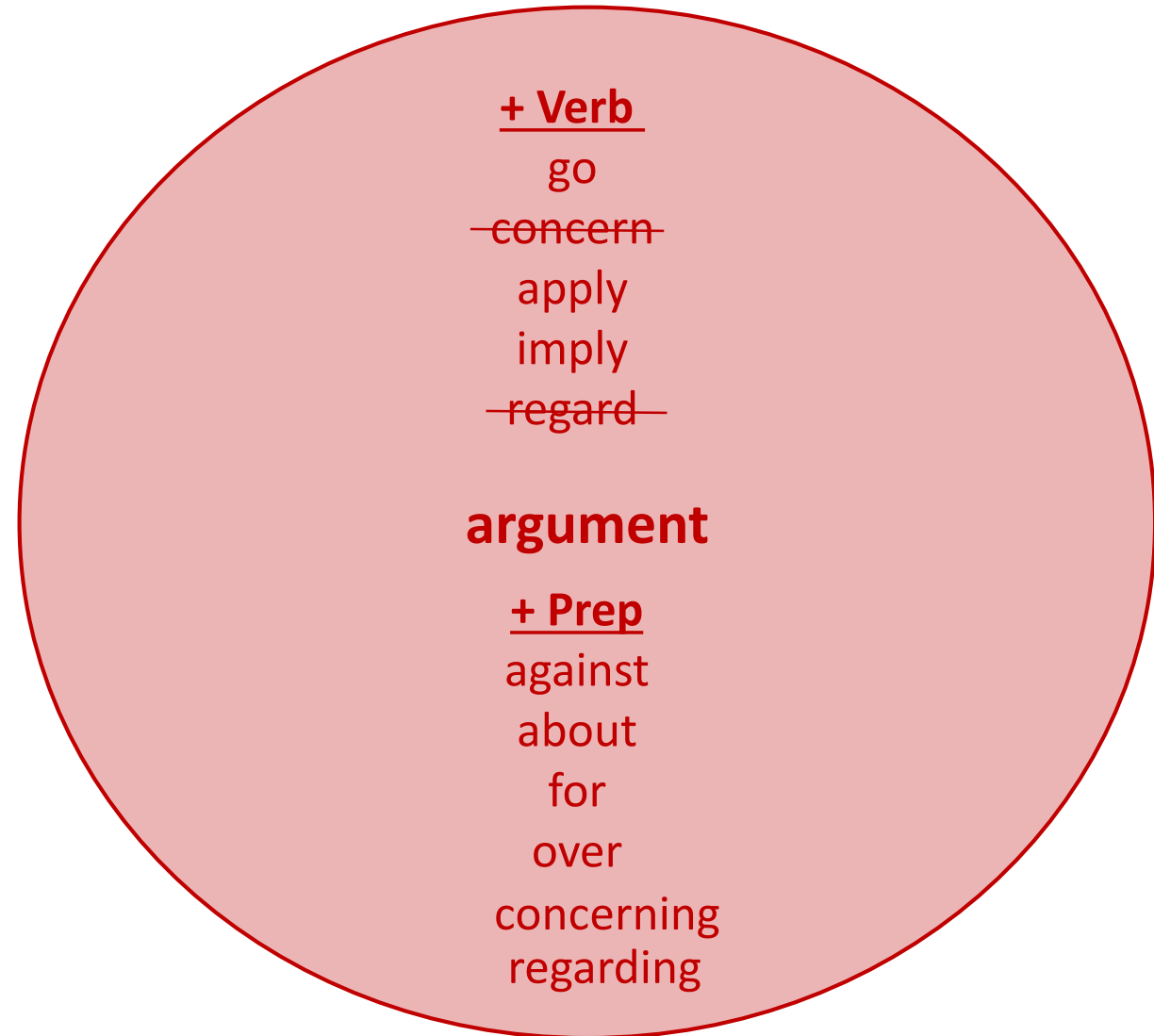
conduct	180	10.56	...
undertake	193	10.42	...
fund	24	8.08	...
publish	33	7.69	...
support	52	7.34	...
pursue	17	7.19	...
cite	12	6.77	...
sponsor	8	6.76	...
stem	8	6.76	...
stimulate	9	6.6	...
summarize	7	6.51	...

Verbs used with *research* as object in PICAE

Collocation research issues

Another problem...

- Word sketches sort collocates according to word class
- But collocation paradigms don't always evoke a single POS
- For example:
 - Their **argument** **about** film and dream*
 - Their **argument** **over** film and dream*
 - Their **argument** **concerning** film and dream*
 - Their **argument** **regarding** film and dream*
- Decision to move them



Collocation research issues

1. attitude + to

- *Positive **attitudes to** technology*
- ~~*Convey the **attitude to** the student that...*~~
- *Convey the attitude **to** the student that...*

2. community + develop

- *How **communities develop** through time*
- ~~*Connectedness between the **community develops***~~
- *Connectedness between the community **develops***

Similar on surface, but grammar relation differs

Collocate misidentified (distorting logDice & co-occurrence figures)

Collocation research issues

1. Late (ADJ) + development

- *the late development of passives*
- *the latest technological developments in the field*

2. Increase (ADJ) + demand

- *A momentarily increased demand for...*
- *To satisfy increasing energy demands ...*

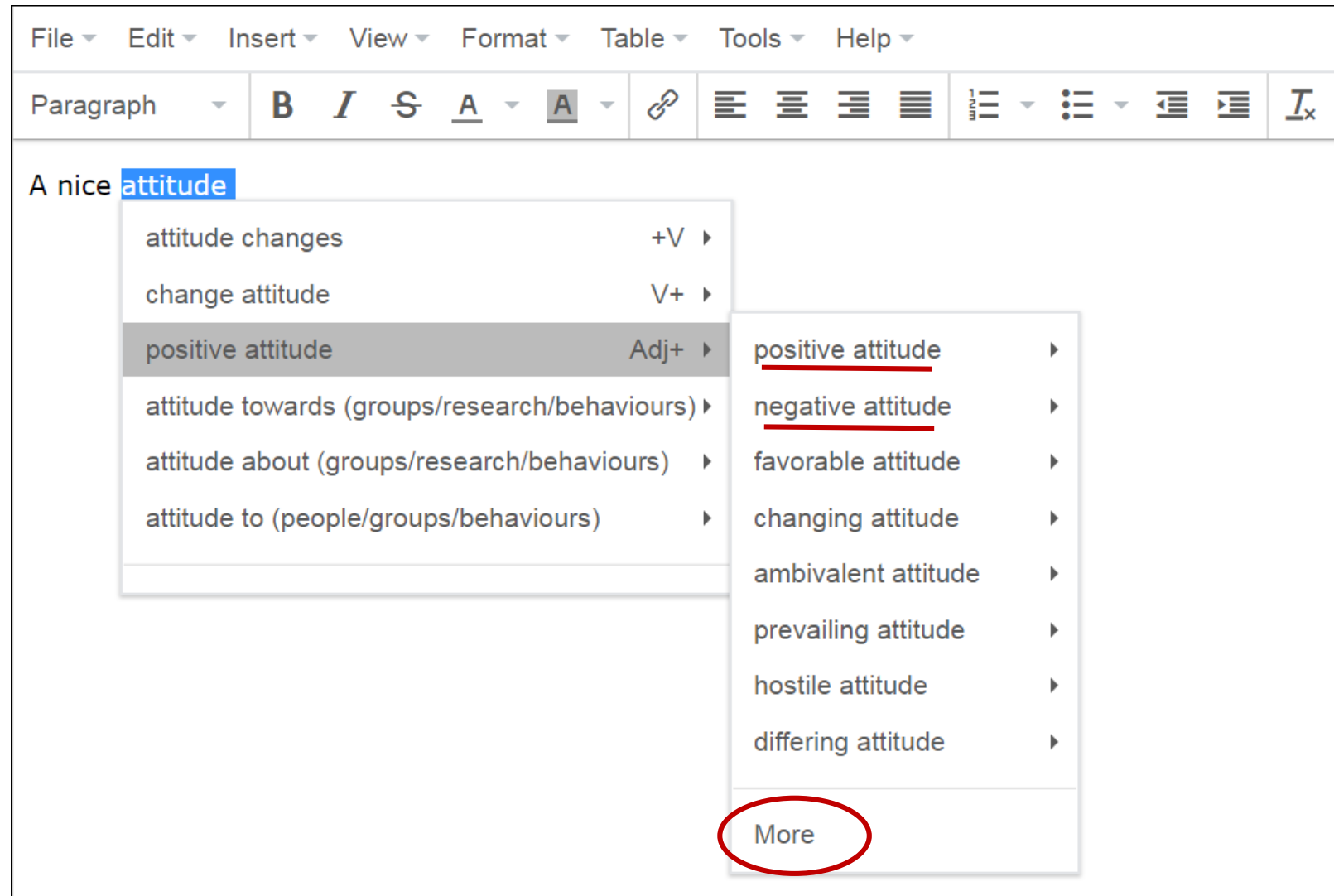
Similar on the surface, but semantically different

- Affects logDice & co-occurrence
- Polysemous collocates get bumped up the queue

Collocation research issues

To avoid overwhelming writers with too many collocations...

At this level, we display only the top eight collocates (sorted by logDice)



The screenshot shows a word processing application interface. The menu bar includes File, Edit, Insert, View, Format, Table, Tools, and Help. The Paragraph menu is open, showing options for Bold (B), Italic (I), Strikethrough (ABC), Underline (A), and another Underline (A) option. The text "A nice attitude" is visible, with "attitude" highlighted in blue. A dropdown menu is open for "attitude", listing the following collocates with their grammatical categories:

- attitude changes +V
- change attitude V+
- positive attitude Adj+ (highlighted)
- attitude towards (groups/research/behaviours)
- attitude about (groups/research/behaviours)
- attitude to (people/groups/behaviours)

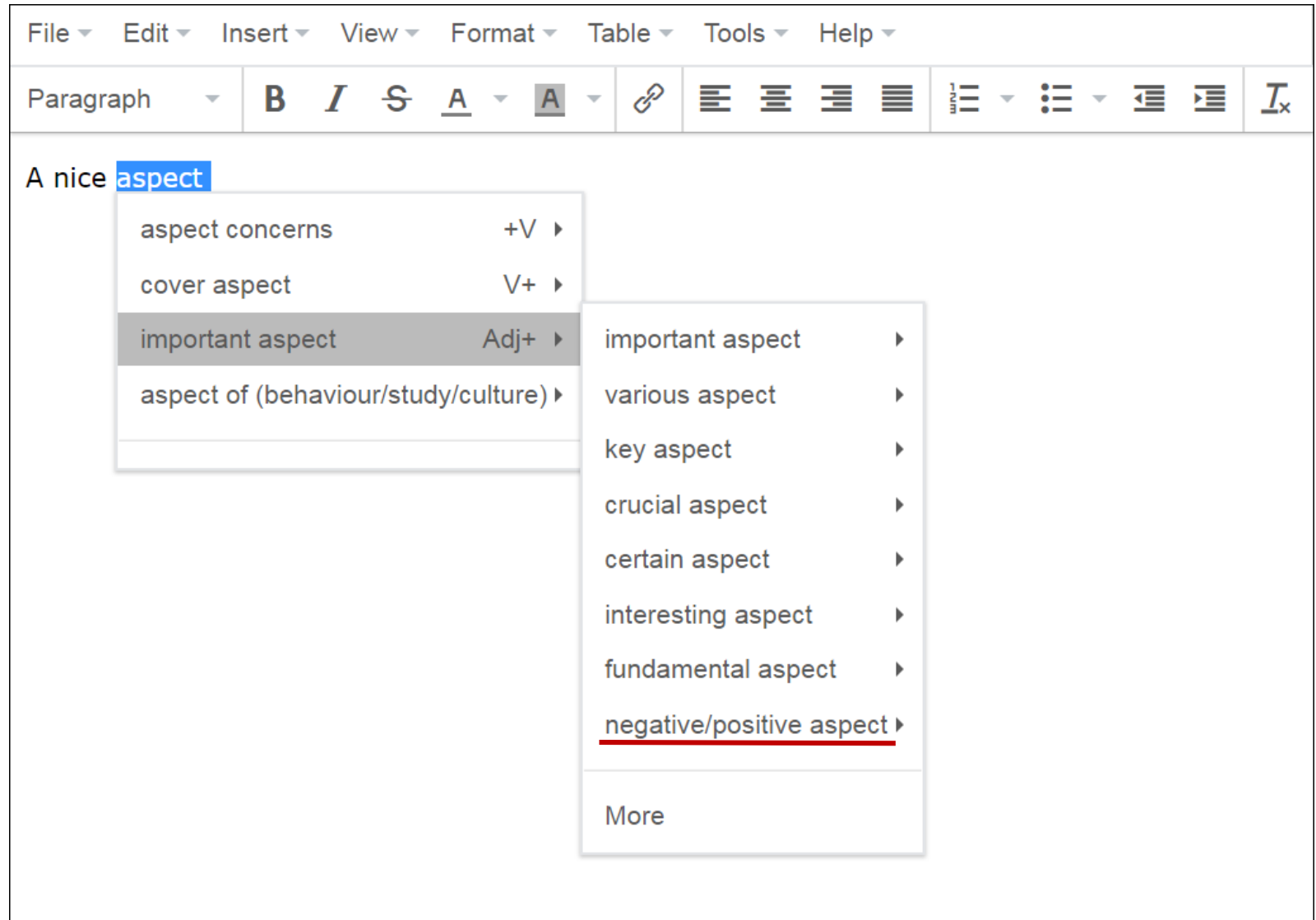
A second dropdown menu is open for "positive attitude", listing the following collocates:

- positive attitude
- negative attitude
- favorable attitude
- changing attitude
- ambivalent attitude
- prevailing attitude
- hostile attitude
- differing attitude

A "More" button is circled in red at the bottom of the second dropdown menu.

Collocation research issues

However, manual curation was necessary in certain cases, to avoid impression of inconsistent coverage



The screenshot shows a word processing application interface with a menu open for the word "aspect". The menu lists several collocations with their grammatical structures:

- aspect concerns +V ▶
- cover aspect V+ ▶
- important aspect Adj+ ▶
- aspect of (behaviour/study/culture) ▶

A secondary menu is open for "important aspect", listing various adjectives:

- important aspect ▶
- various aspect ▶
- key aspect ▶
- crucial aspect ▶
- certain aspect ▶
- interesting aspect ▶
- fundamental aspect ▶
- negative/positive aspect ▶

The "negative/positive aspect" option is underlined in red. A "More" button is visible at the bottom of the secondary menu.

Example selection

Example selection

1. Examples selected from corpora of expert academic writing
 - From word sketch to concordance
 - With editorial intervention (to select, adapt, anonymize, shorten)
2. Are interdisciplinary or transferable to other disciplines
3. Promote data-driven learning (Johns 1991)
 - Three examples (Frankenberg-Garcia 2014)
 - Colligation cues if relevant/possible
4. Brief: less effort to process, less distracting for writers
 - GDEX (Kilgarriff et al. 2008)

Example selection issues

V obj N*	
undertake	...
activities undertaken	
coordinate	...
regulate	...
perform	...
organize	...
inhibit	...
monitor	...
increase	...
stimulate	...



<s> But by a clever construction of notional prices (called ' shadow prices ' ; Chapters 7-8) , economists have adapted GDP even for economies like Desta 's , where much economic **activity** is **undertaken** in non-market institutions . </s>

<s> Because people can't insure themselves sufficiently against failure , they are reluctant to **undertake activities** offering a chance of huge success if there is also an accompanying chance of large failure . </s>

<s> Once all of the goals and indicators for success (for the collaboration , for the school , for the community , for the parents , and for the students) have been identified and it is clear what has to happen first , interventions (**activities undertaken** in the community school , such as after-school programs , health services , or parent workshops) are added wherever needed to reach an outcome . </s>

<s> In general , online users showed a greater predilection to expand the set of online **activities undertaken** in a year 's time . </s>

<s> This is important to lifelong learning where self and peer assessment provide critical and informal feedback in the workplace and in the myriad of other human **activities undertaken** in daily life . </s>

Example selection issues

V obj N*	
undertake	...
activities undertaken	

GDEX →

Longest commonest match (colligation) can be less evident

perform	...
organize	...
inhibit	...
monitor	...
increase	...
stimulate	...

<s> Companies **undertake** marketing **activities** in order to elicit some kind of response from buyers . </s>

<s> Having difficulty **undertaking** a particular **activity** is inherently a subjective determination . </s>

<s> Rather , there has been the specialization of the **activities undertaken** by business units and increasing cooperation between them . </s>

<s> Edison 's organizational innovation lay in the range and scale of the research **activities undertaken** . </s>

<s> They **undertake** political **activity** , but not of the traditional kind . </s>

<s> For patients at high risk of a cardiovascular event , **undertaking** sexual **activity** may constitute a significant risk . </s>

<s> When an international firm crosses national borders to **undertake** business **activities** , it may also cross cultural borders . </s>

<s> Applicants must only **undertake** unpaid **activities** related to the purposes of the sponsoring organization . </s>

ColloCaid 0.4

Version 0.4

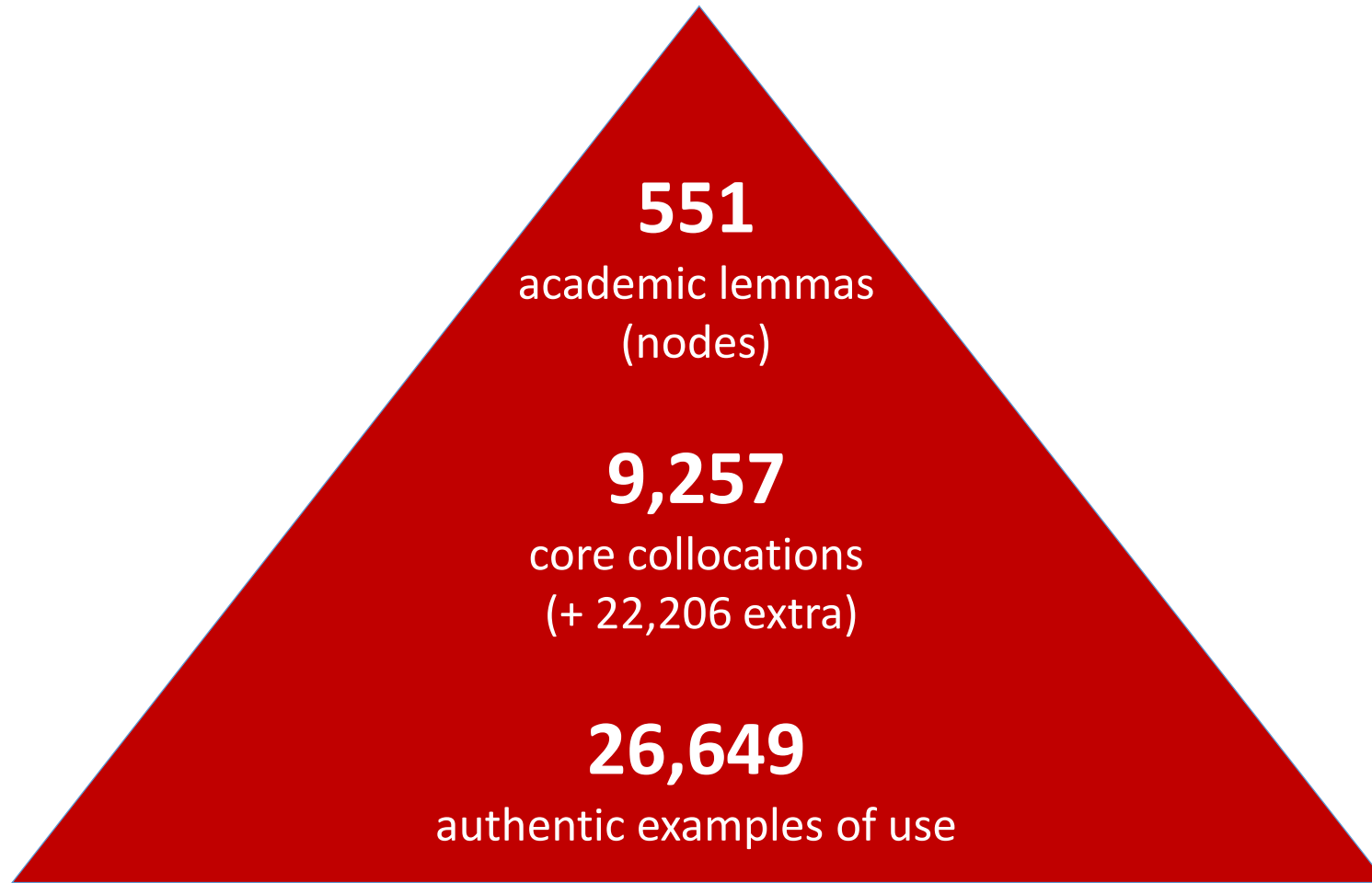
File ▾ Edit ▾ Insert ▾ View ▾ Format ▾ Table ▾ Tools ▾ Help ▾

Paragraph ▾ **B** *I* ~~S~~ A ▾ A ▾

Another **activity**

- activity generates +V ▾
- undertake activity V+ ▾**
 - activities undertaken ▾**
 - an ability to function and **undertake** usual activities
 - airports **undertake** public relations activities via websites
 - children can **undertake** such developmental activities
 - ** Please take care to adapt examples to your own text **
 - coordinate activities ▾
 - regulate activities ▾
 - perform activities ▾
 - organize activities ▾
 - monitor activity ▾
 - increase activity ▾
 - control activity ▾
 - More
- increased activity Adj+ ▾
- activity of (organisation/body part/group) ▾
- activity within (place/group) ▾
- activity during (stage/period) ▾
- activity among (groups/individuals) ▾
- activity outside (home/family) ▾
- activity in (body part/sector/area) ▾

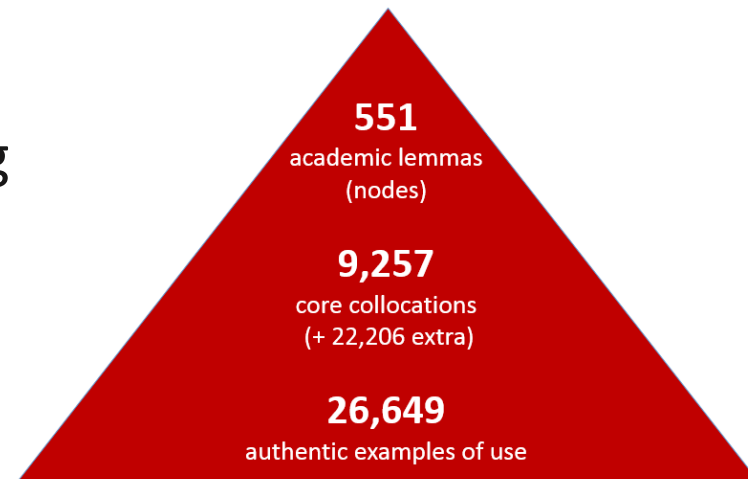
! Database currently being edited and proofread



Conclusion

Conclusion

- We couldn't have achieved anything near what we have managed so far without e-lexicographic tools & resources
- But spent a lot of time curating the data manually
- We hope our analysis of what needs editing/manual curation can contribute to further advances in e-lexicography
 - Collocation Workshop 30/09/2019 😊
 - e.g. Rychlý & Jakubíček; Fuhrmann et al.; Krek et al.
- And we draw attention to the dangers of uncritical use of existing tools and resources
 - Trust the text (Sinclair 2004), but beware of the rest



The screenshot shows the Collocaid web application interface. At the top left is the 'collocaid' logo, and at the top right is the user name 'Ana' with a dropdown arrow. Below the header is a menu bar with options: File, Edit, Insert, View, Format, Table, Tools, and Help. Underneath the menu bar is a toolbar with icons for Paragraph, Bold (B), Italic (I), Strikethrough (S), Underline (A), Background Color (A), Link, Text Alignment (left, center, right, justified), Bulleted List, Numbered List, Decrease Indent, Increase Indent, and Text Color (I_x).

The main text area contains the sentence: "Thank you for your **attention**". The word "attention" is highlighted in blue. A dropdown menu is open below the word, listing various collocations:

- attention shifts +V ▶
- pay attention V+ ▶
- little attention Adj+ ▶
- little attention ▶
- much attention ▶
- special attention ▶
- close attention ▶
- careful attention ▶
- considerable attention ▶
- scholarly attention ▶
- particular attention ▶
- More

The "careful attention" option is highlighted in grey. A secondary dropdown menu is open to the right of "careful attention", showing example sentences:

- they give each and every manuscript **careful attention**
- the manager must give this matter **careful attention**
- the consequences of this system deserve **careful attention**

At the bottom of this secondary menu, there is a note: "** Please take care to adapt examples to your own text **".

Further information

www.collocaid.uk

AHRC - AH/P003508/1

Principal Investigator Dr Ana Frankenberg-Garcia (Surrey)

Co-investigators Prof Robert Lew (Poznan), Prof Jonathan Roberts (Bangor)

Researchers Dr Geraint Rees (Surrey), Dr Nirwan Sharma (Bangor)

& Peter Butcher (Bangor)

Come to our demo tomorrow!