**Corpus Tools Brainstorming Session**
What tools do we need for the future?

**Laurence ANTHONY**
Professor, Center for English Language Education,
Waseda University, Japan
Honorary Research Fellow (2013-2015), Lancaster University

anthony@waseda.jp
http://www.antlab.sci.waseda.ac.jp/software/

AACL 2014, Northern Arizona University, Flagstaff, US
September 25-28, 2014

Faculty of Science and Engineering, Waseda University

---

## Overview of this brainstorming session

- Background
  - Reasons for organizing this session
  - Possible outputs from this session
  - The current state of corpus linguistics tools
  - The case for programming your own corpus tools
- Brainstorming:
  - Part 1: What do you currently do with corpora
  - Part 2: What do you want to do with corpora (that you cannot already)?
  - Part 3: What prevents you collaborating with a tools developer to create new tools… or programming your own?
  - Part 4: How can we increase corpus linguists' attention to corpus tools?

2

---

## Background:
Reasons for organizing this session

3

---

## Background:
Reasons for organizing this session

- To highlight the importance of tools in our field
  - Many corpus linguists forget that they look through the lens of a corpus tool when doing their research
- To introduce and discuss the current solutions to corpus linguists' needs
  - Many corpus linguists are not aware of currently available tools
- To generate ideas for new (and useful) corpus tools
  - What functions of a corpus tool are important to you?
  - What would you like corpus tools to do?
- But…
  - This is not a session to voice criticisms of current tools
    - (Unless the criticism is of *AntConc*!)

4

---

## Background:
Possible outputs from this session

- The development of new corpus tools
  - e.g. realization of ideas generated in this session by current developers
  - e.g. collaborative tools development between developers and session participants
  - e.g. creation of new tools developers with an interest in programming
- A regular corpus tools section in corpus linguistics conferences
- More research presentations/publications on corpus tool innovations

5

---

## Background:
The current state of corpus linguistics tools

6

1

## Background:
### The current state of corpus linguistics tools

- A definition of corpus linguistics
    - It is an empirical (experimental) approach
        - an analysis of actual patterns of use in target texts
    - It uses a corpus of natural texts as the basis for analysis
        - a representative sample of target language stored as an electronic database
    - It relies on computer software for analysis
        - results are generated using automatic and interactive techniques
    - It depends on both quantitative and qualitative analytical techniques
        - observations are counted and results are interpreted

7

(Biber et al., 1998)

---

## Background:
### The current state of corpus linguistics tools

- Four Generations of Corpus Tools
  (see McEnery & Hardie, 2012)
    - 1st-generation (1960s-1970s)
        - run on mainframes, ASCII-based, very limited functions
            - e.g., *A Concordance Generator (Smith, 1966)*
            - *e.g., Discon (Clark, 1966)*
            - *e.g., Drexel Concordance Program (Price, 1966)*
            - *e.g., Concordance (Dearing, 1966)*
            - *e.g., CLOC* (Reed, 1978)

8

---

## Background:
### The current state of corpus linguistics tools

- *Discon (Clark, 1966)*

*Discon.* Purpose of program: Concordance-making. This program is simply the well-known DISCON, originally written for the 7090 and converted to the 7040 by R. L. Priore of Roswell Park Memorial Institute, Buffalo, New York, then converted to the 7044 here by Marjorie Schultz. Type and format of input: Punched cards: 6 cols. ID, the rest data.
Programming language used: Fortran IV. Required hardware: IBM 7044. Running time: Approx. 4 min. per 1000 lines of poetry, exclusive of printout time.
Correspond with Roger Clark or Lewis Sawin, 123 W Hellems, University of Colorado, Boulder, Colorado.

Computers and the Humanities (1966, Vol. 1, Issue 2, p. 39)

9

---

## Background:
### The current state of corpus linguistics tools



10

---

## Background:
### The current state of corpus linguistics tools

- Four Generations of Corpus Tools
    - 2nd-generation (1980s-1990s)
        - run on PCs, ASCII-based, limited functions, scalability problems
            - e.g., *Oxford Concordance Program (OCP)* (Hockey, 1988)
            - e.g., *Longman Mini-Concordancer* (Chandler, 1989)
            - e.g., *Kaye concordancer* (Kaye, 1990)
            - *e.g., MicroConcord* (Scott & Johns, 1993)



MicroConcord (Scott & Johns, 1993)

11

---

## Background:
### The current state of corpus linguistics tools



Tim Johns and Randolph Quirk (1982/1983)
Photo by John Higgins.

12

---

2

## Slide 13

# Background:
## The current state of corpus linguistics tools

- Four Generations of Corpus Tools
  - 3rd-generation (2000s-present)
    - more functions, better statistics, improved scalability, multi-language support , more user-friendly, simple, flexible
      - e.g., *WordSmith Tools* (Scott, 1996-2014)
      - e.g., *MonoConc Pro* (Barlow, 2000)
      - e.g., *AntConc* (Anthony, 2004-2014)
    - Problems:
      - limited functionality (still)
      - limited access to corpora containing copyrighted data
      - limited ability to scale to massive (100 million+) corpora

13

## Slide 14

# Background:
## The current state of corpus linguistics tools



WordSmith Tools (Scott, M. , 2014)

14

## Slide 15

# Background:
## The current state of corpus linguistics tools



MonoConc Pro (Barlow, M. , 2000)

15

## Slide 16

# Background:
## The current state of corpus linguistics tools



AntConc (Anthony, L., 2014)

16

## Slide 17

# Background:
## The current state of corpus linguistics tools

- Four Generations of Corpus Tools
  - 4th-generation (late 2000s-present)
    - better scalability, access to copyrighted data, available anywhere or everywhere (via browsers)
      - e.g., corpus.byu.edu (Davies, 2011), *CQPweb* (Hardie, 2011), *SketchEngine* (Kilgariff, 2011), *Wmatrix* (Rayson, 2011)
    - Problems:
      - limited functionality (still)
      - overkill for many purposes (e.g., analyzing small corpora)
      - logistics: registration, setup, licenses, payment
      - limited access to public-domain corpora
      - no access to "personalized" corpora
        - e.g., institution-owned collections of copyrighted material
        - e.g., user-created collections of copyrighted material
      - explosion in one-off, web-based, single-corpus interfaces
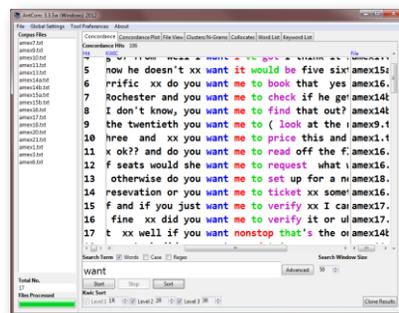
17

## Slide 18

# Background:
## The current state of corpus linguistics tools



18

3

## Background:
### The current state of corpus linguistics tools



19

## Background:
### The current state of corpus linguistics tools



| | |
|---|---|
| corpus.byu.edu | |
| AntConc | |
| WordSmith Tools | |
| Sketch Engine | |
| Sarah (with BNC) | |
| Monoconc Pro | |
| Xaira (with BNC XML or your own... | |
| WMatrix | |
| Oxford Concordancing Program | |
| Longman Mini-concordance | |
| Other | |

"Which computer programs do you use for analysing corpora?"
International survey of corpus linguists. Reponses: 891. (Tribble, 2012)

20

## Background:
### The case for programming your own corpus tools

21

## Background:
### The case for programming your own corpus tools

- Corpus linguists should learn to program ... (?)
    - (Biber, Gries, Weisser, ....)



22

## Background:
### The case for programming your own corpus tools

If you program ...

"you can do analyses not possible with concordancers ...
you can do analyses more quickly and more accurately ...
you can tailor the output to fit your own research needs ...
you can analyze a corpus of any size"

(Biber et al., 1998, p. 256)

23

## Background:
### The case for programming your own corpus tools

"when you use pre-configured corpus programs, you're a little bit at the mercy of the company or individual selling them ...

One final big advantage of programming languages, therefore, is that you are in the driver's seat."

(Gries, 2009, p. 11-12)

24

4

## Background:
### The case for programming your own corpus tools

The reality for most corpus researchers, however, is that computer programming is in a completely different world ... without extensive training in programming ... it is likely that these [DIY] tools would be more restrictive, slower, less accurate and only work with small corpora.

(Anthony, 2009, p. 95)

25

## Background:
### The case for programming your own corpus tools

But, always remember ...

"Research should be led by the science not the tool."

Professor Jim Wild, Lancaster University
Vice-President, Royal Astronomical Society

26

## Brainstorming Part 1:
What do you currently do with corpora?

27

## Brainstorming Part 1:
### What do you currently do with corpora?

- Write down everything you do with corpora now?
  - Try to think about what you DO with corpora (not what tool you use)
    - e.g. identify common word/phrase patterns (in context)
    - e.g. find unusually frequent words/phrases in the corpus (i.e. keywords)
  - Think about WHY you do these things with corpora
    - e.g. to help EFL students use standard English
    - e.g. to identify characteristic features of a text/genre

28

## Brainstorming Part 1:
### What do you currently do with corpora?

- Participant Responses (added after discussion)

| find word/phrase patterns (KWIC) | match patterns in text (via scripting) |
| find word/phrase positions (Plot) | generate statistics (e.g. using R) |
| find collocates | measure dispersion of word/phrase patterns |
| find N-grams/Lexical bundles | compare words/synonyms |
| find Clusters | identify characteristics of texts |
| generate word lists | |
| generate keyword lists | |

29

## Brainstorming Part 2:
What do you want to do with corpora
(that you cannot already)?

30

5

## Brainstorming Part 2:
### What do you want to do with corpora?

- Write down interesting things you want to do with a corpus?
  - Do not worry if the idea is crazy or impossible.
  - But, also consider WHY you want to do this.
    - e.g. Find topic sentences in paragraphs
      - WHY – As a source of examples for writing classes

31

## Brainstorming Part 2:
### What do you want to do with corpora?

- Participant Responses (added after discussion)

| | |
|---|---|
| compute distances between subsequent occurrences of search patterns (words, lemmas, POS, ...) | process audio data |
| quantify the degree of variability around search patterns | carry out phonological analysis (e.g. neighbor density) |
| generate counts per text (in addition to corpus) | use tools to build a corpus (e.g. finding texts, annotating texts, converting non-ascii characters to ascii) |
| extract definitions | create new visualizations of data (e.g. a roman candle of words that 'explode' out of a text) |
| find patterns of range and frequency | identify the encoding of corpus texts |
| work with private data but allow for powerful handling of annotation (e.g. comparing frequencies of sub-corpora) | compare two corpora along many dimensions |
| carry out extensive move analysis over large texts | identify changes in language over time |
| search corpora by semantic class | disambiguate word senses |

32

## Brainstorming Part 3:
What prevents you collaborating with a tools developer to create new tools…
or programming your own?

33

## Brainstorming Part 3:
### What prevents you collaborating with a tools developer to create new tools… or programming your own?

- Participant Responses (added after discussion)

| Collaboration | Programming |
|---|---|
| not confident to contact developers directly | not enough time to learn programming skills |
| not sure what features are already available in current tools | the need for programming is not immediately apparent |
| | powerful tools already exist |
| | people are happy to do what they *can* with corpora instead of doing what they *should* do with corpora |

34

## Brainstorming Part 4:
How can we increase corpus linguists' attention to corpus tools?

35

## Brainstorming Part 4:
### How can we increase corpus linguists' attention to corpus tools?

- Participant Responses (added after discussion)

| | |
|---|---|
| continue holding brainstorming sessions of this kind | introduce corpus tools tracks in conference programs |

36