# La linguistique de corpus à l'appui des synthèses de recherche

alex.boulton@univ-lorraine.fr

ATILF, 22 mars 2024

Journée thématique et transversale :
 *Linguistique de corpus à la croisée de questionnements théoriques, méthodologiques et empiriques*
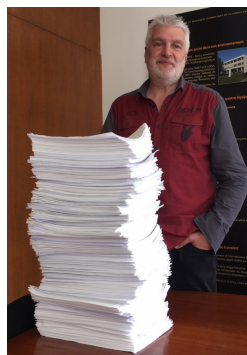
aka 'data-driven learning' (DDL)

"the attempt to cut out the middleman as far as possible and...
give the learner direct access to the data" (Johns, 1990, p.18)

"using the tools and techniques of corpus linguistics
for pedagogical purposes" (Gilquin & Granger, 2010, p.359)     ⇨ ±directly for L2

My empirical DDL collection:

2007 =   39

2012 = 116

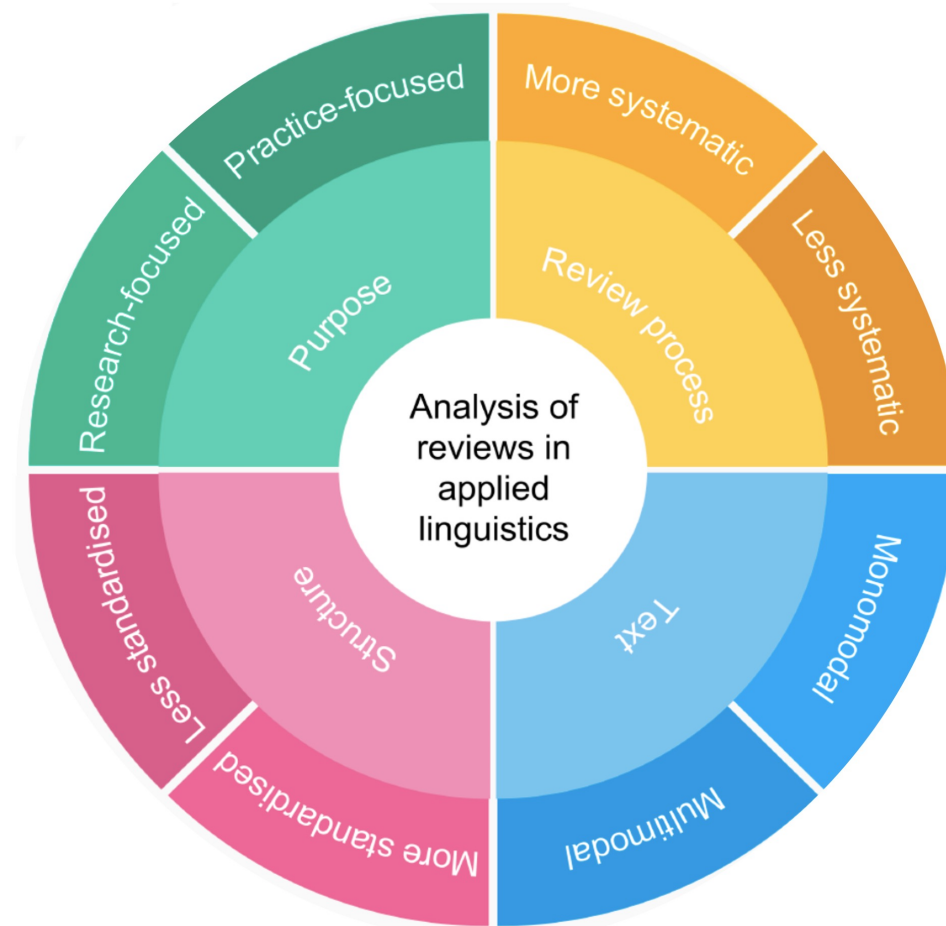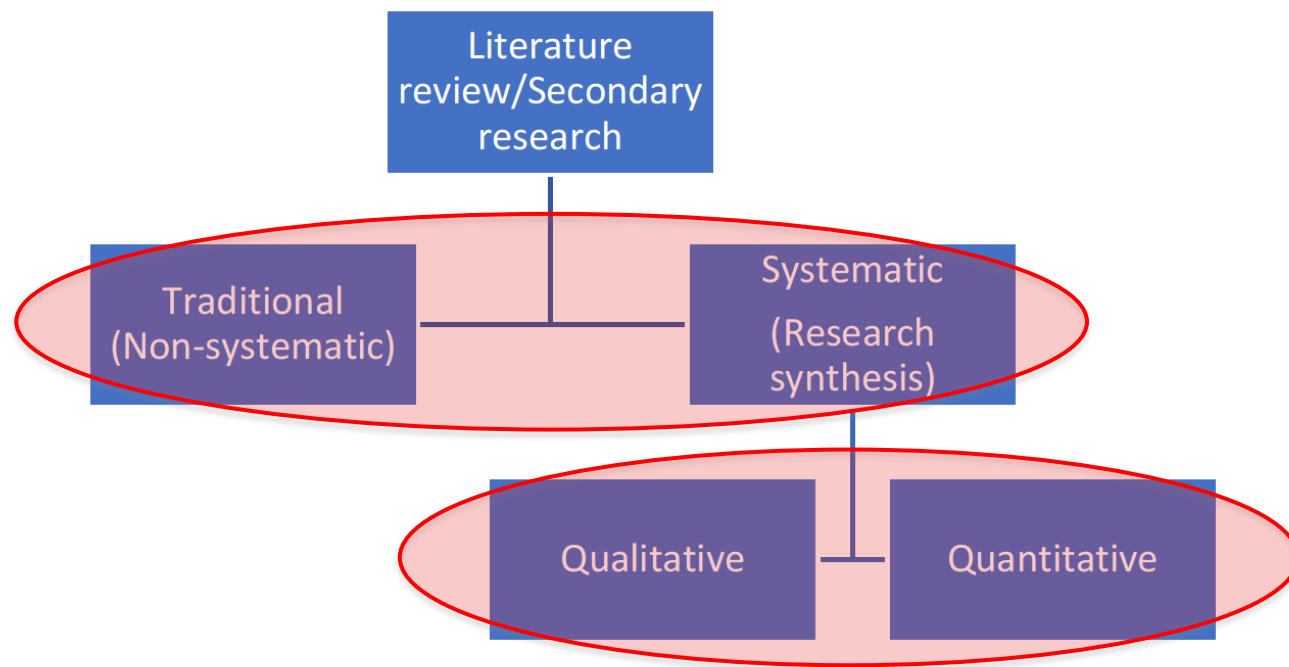2017 = 210

2019 = 351

2021 = 489...

**now**
**777!**

syntheses...

## A typology of secondary research in applied linguistics
(Chong & Plonsky, advance access)



1. critical review; 2. meta-analysis; 3. methodological synthesis; 4. mixed review; 5. narrative review;
6. qualitative research synthesis; 7. research agenda; 8. research into practice; 9. scoping review;
10. state-of-the-art review; 11. systematic literature review; 12. historical review; 13. bibliometric review

## Qualitative (narrative)

- 2007 Chambers (12 studies)
- 2007 Boulton (39 studies)
- 2010 Boulton (27 studies, learning outcomes)
- 2011 Yoon (12 studies, concordancing)
- 2012 Boulton (20 studies, ESP)
- **2013 Boulton & Tyne (116 studies)**
- 2017 Luo & Zhou (18 studies, writing)
- 2017 Boulton (46 studies, historical timeline)
- 2018 Chen & Flowerdew (37 studies, EAP)
- 2019 Al-Gamal & Ali (5 studies, recent)
- 2023 Sun & Park (32, collocations)

## Quantitative (meta-analyses)

- 2015 Mizumoto & Chujo (14 studies, Japan)
- 2015 Cobb & Boulton (21 studies, preliminary)
- 2017 Boulton & Cobb (64 studies)
- 2019 Lee et al. (29 studies, vocab)
- 2023 Ueno & Takeuchi (144 studies)

## Other (mixed)

- 2019 He & Wei (328 studies, bibliometric)
- 2021 Boulton (351 studies, coding)
- 2021 Boulton & Vyatkina (489 studies, scoping)
- 2022 Pérez-Paredes (32 studies, keywords/clusters)
- 2023 Dong et al. (126 studies, bibliometric)
- 2023 Lusta et al. (89 studies, systematic review)
- 2024 Boulton & Vyatkina (148 studies, English, JIF)

☺ Wide-ranging, rich, in-depth

☹ Cherry-picking, subjective

**Qualitative (narrative)**
- 2007 Chambers (12 studies)
- 2007 Boulton (39 studies)
- 2010 Boulton (27 studies, learning outcomes)
- 2011 Yoon (12 studies, concordancing)
- 2012 Boulton (20 studies, ESP)
- 2013 Boulton & Tyne (116 studies)
- 2017 Luo & Zhou (18 studies, writing)
- 2017 Boulton (46 studies, historical timeline)
- 2018 Chen & Flowerdew (37 studies, EAP)
- 2019 Al-Gamal & Ali (5 studies, recent)
- 2023 Sun & Park (32, collocations)

**Quantitative (meta-analyses)**
- 2015 Mizumoto & Chujo (14 studies, Japan)
- 2015 Cobb & Boulton (21 studies, preliminary)
- **2017 Boulton & Cobb (64 studies)**
- 2019 Lee et al. (29 studies, vocab)
- 2023 Ueno & Takeuchi (144 studies)

**Other (mixed)**
- 2019 He & Wei (328 studies, bibliometric)
- 2021 Boulton (351 studies, coding)
- 2021 Boulton & Vyatkina (489 studies, scoping)
- 2022 Pérez-Paredes (32 studies, keywords/clusters)
- 2023 Dong et al. (126 studies, bibliometric)
- 2023 Lusta et al. (89 studies, systematic review)
- 2024 Boulton & Vyatkina (148 studies, English, JIF)

**EMPIRICAL STUDY**

# Corpus Use in Language Learning: A Meta-Analysis

Alex Boulton and Tom Cobb

Université de Lorraine and Université du Québec à Montréal

**LANGUAGE LEARNING**

To see:
a) *if* DDL works
b) *how well* DDL works
c) *where* DDL works (...or doesn't)

☺Quantitative: rigorous, pooled data for clear answers
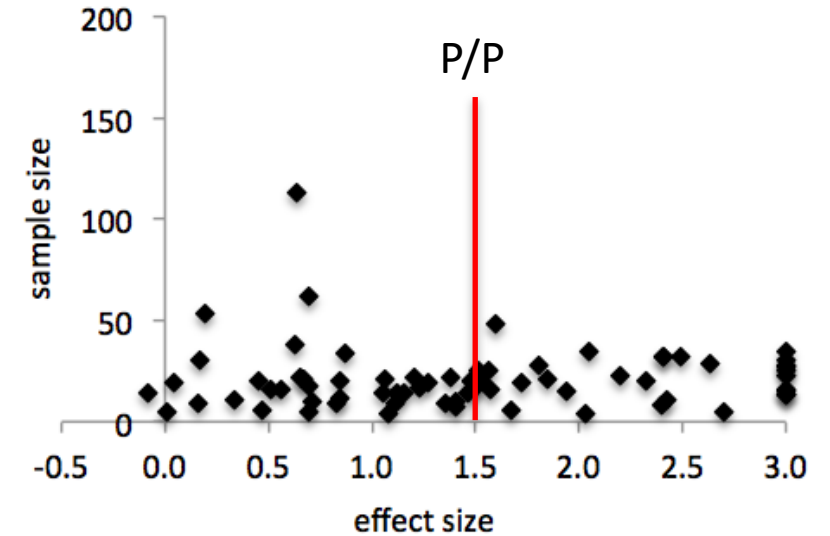☹Quantitative: less inclusive, less nuanced, subjective

$$d = \frac{M_2 - M_1}{\sqrt{\frac{SD_1{}^2 + SD_2{}^2}{2}}}$$

| Effect size | Plonsky & Oswald 2014 (C/E, *n*=67) | Plonsky & Oswald 2014 (P/P, *n*=25) | cf. SLA |
|---|---|---|---|
| large | 0.9 | 1.4 | 1st quartile |
| medium | 0.6 | 1.0 | 2nd quartile |
| small | 0.4 | 0.6 | 3rd quartile |

| Boulton & Cobb 2017 | 0.95 (*k*=50) | 1.50 (*k*=71) |
|---|---|---|

☺☺☺☺☺☺☺☺

DDL large effects. DDL good.
End of story. Everyone go home.



Moderator Variables:
  "DDL works pretty well
  in almost any context
  where it has been
  extensively tried." (p. 386)
But…

Qualitative (narrative)
- 2007 Chambers (12 studies)
- 2007 Boulton (39 studies)
- 2010 Boulton (27 studies, learning outcomes)
- 2011 Yoon (12 studies, concordancing)
- 2012 Boulton (20 studies, ESP)
- 2013 Boulton & Tyne (116 studies)
- 2017 Luo & Zhou (18 studies, writing)
- 2017 Boulton (46 studies, historical timeline)
- 2018 Chen & Flowerdew (37 studies, EAP)
- 2019 Al-Gamal & Ali (5 studies, recent)
- 2023 Sun & Park (32, collocations)

Quantitative (meta-analyses)
- 2015 Mizumoto & Chujo (14 studies, Japan)
- 2015 Cobb & Boulton (21 studies, preliminary)
- 2017 Boulton & Cobb (64 studies)
- 2019 Lee et al. (29 studies, vocab)
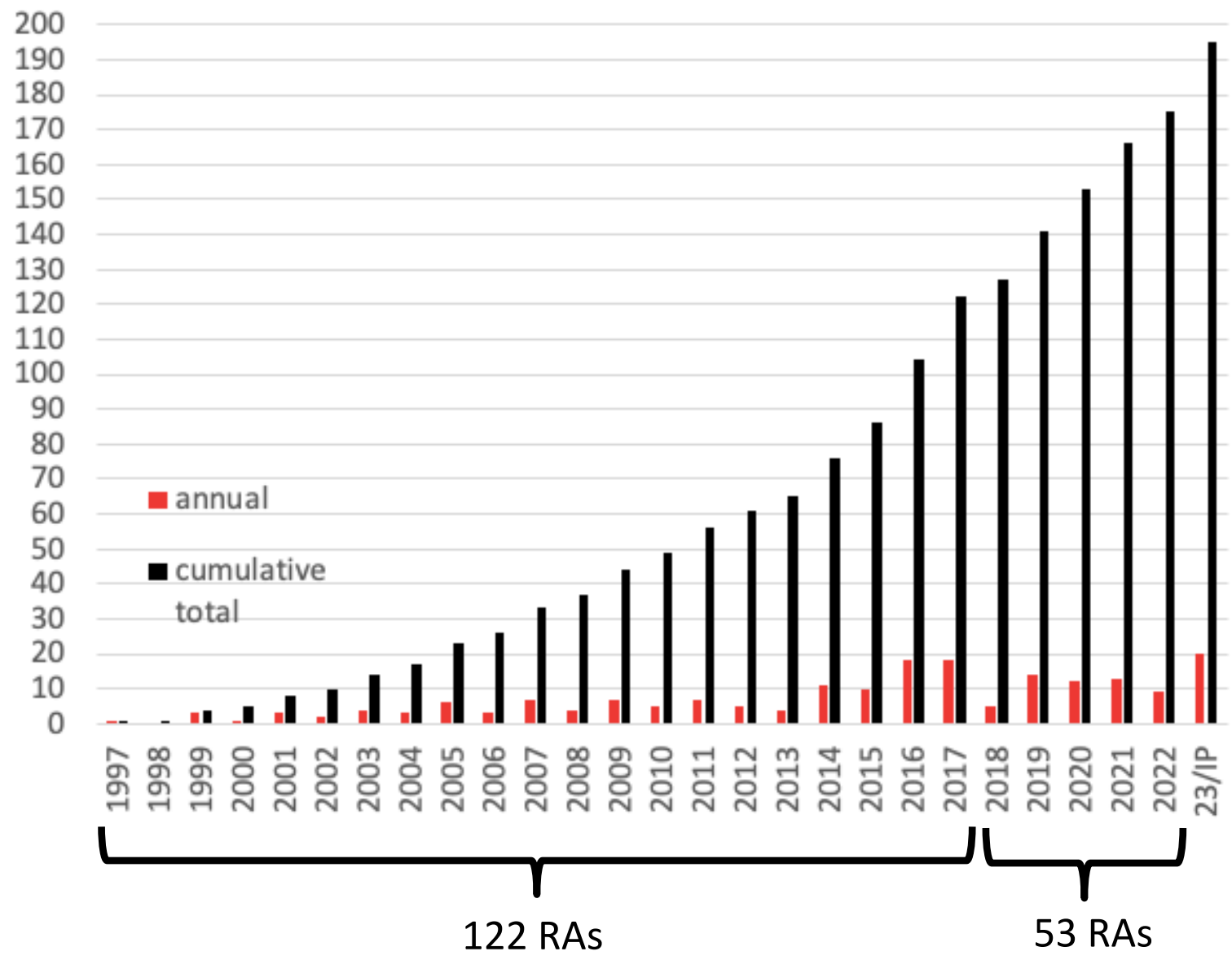- 2023 Ueno & Takeuchi (144 studies)

Other (mixed)
- 2019 He & Wei (328 studies, bibliometric)
- 2021 Boulton (351 studies, coding)
- **2021 Boulton & Vyatkina (489 studies, scoping)**
- 2022 Pérez-Paredes (32 studies, keywords/clusters)
- 2023 Dong et al. (126 studies, bibliometric)
- 2023 Lusta et al. (89 studies, systematic review)
- **2024 Boulton & Vyatkina (148 studies, English, JIF)**

- Methodical collection of published empirical DDL studies
  (cf. Boulton & Cobb, 2017; Boulton, 2021; Boulton & Vyatkina, 2021; Boulton & Vyatkina, 2024)

- Today: up to 2022 inclusive (thanks to A. Jakob Johnson)
  DDL, empirical, in English, JCR-ranked LING+EDU  journals
  ☺ ±exhaustive, but… ☹ what's NOT included
  ☺ highly visible, but…  ☹ impact factor ≠ not quality!  Clarivate Analytics

In the last 5 years (2018-2022):
  RQ1. What trends are emerging in DDL research?
        coding and analysis – manual
  RQ2. How do researchers talk about DDL?
        corpus analysis ('aboutness') – AntConc

# RA corpus & timeline



| Title | 1997-2017 | 2018-2022 | TOTAL |
|---|---|---|---|
| ReCALL | 26 | 6 | 32 |
| CALL | 24 | 3 | 27 |
| LLT | 23 | 2 | 25 |
| System | 8 | 6 | 14 |
| JEAP | 3 | 6 | 9 |
| IJAL | 3 | 5 | 8 |
| ESP | 4 | 3 | 7 |
| ELTJ | 5 | 1 | 6 |
| IJLex | 4 | 2 | 6 |
| JSLW | 3 | 2 | 5 |
| JCAL | 3 | | 3 |
| Lawareness | 3 | | 3 |
| BJET | 1 | 1 | 2 |
| EIT | 2 | | 2 |
| ETS | 2 | | 2 |
| ILE | 2 | | 2 |
| JCHE | 2 | | 2 |
| LTR | 2 | | 2 |
| MLJ | 1 | 1 | 2 |
| Perspective | 1 | 1 | 2 |
| RELC Journal | 2 | | 2 |
| Misc. | 12 | | 12 |
| | **122** | **53** | **175** |

Coding sheet

(cf. B&V 2021 in IRIS repository)

Excerpt JIF 2018-2022

IRR: decisions, decisions…



- **Publication:** ID, reference, abstract, date, JIF, source, tokens
- **Population:** L1, FL/SL, L2, country, region, proficiency, institution, speciality, discipline, LGP…
- **Treatment:** duration, corpora, size (hands-on), software, interaction, item/skills
- **Research design**: sample, instruments, objective (L/R/A/B), data (Q/Q)

Size (main corpus, hands on):
- <1m                31% ⇨   4%
- 1<99m              36% ⇨ 28%
- >100m              32% ⇨ 68%

Variety (hands-on only) today:
- 0 graded, news, literary, textbooks, parallel
- 1 multimodal

Skills (identifiable, multiple):
- writing            56% ⇨ 88%
- reading            16% ⇨ 12%
- speaking            5% ⇨ 20%
- listening           2% ⇨   0%
- translation        21% ⇨   4%

Language focus (identifiable, multiple):
- vocabulary         24% ⇨ 27%
- lexicogrammar      34% ⇨ 37%
- grammar            16% ⇨ 12%
- discourse          10% ⇨ 10%
- correction         15% ⇨ 15%

# RQ2. Corpus

v4.2.4 + v3.5.8 (Anthony, 2023, 2019)
https://www.laurenceanthony.net/software/antconc

See also:
Jablonkai, R.R., Kim, J., & Yan, R. (in press). A corpus approach to systematic literature reviews. In K. Sadeghi (Ed.), *Routledge handbook of technological advances in researching language learning*.

- 175 RAs ⇨ AntFileConverter ⇨ txt (UTF8)
- main text (meta-data, headers/footers, figs/tables, extracts, foot/endnotes, references, appendices, acknowledgements, etc.)
- check! (hyphens, ligatures; X errors)

| | 1997-2017 | 2018-2022 | TOTAL |
|---|---|---|---|
| papers | 122 | 53 | 175 |
| tokens | 778,020 | 359,692 | 1,137,712 |

# Corpus analysis: wordlist

# Corpus analysis: +stoplist

more revealing?

still a lot in common

# Corpus analysis: keyword list

min freq = 5, min range = 5
(both corpora)

# Corpus analysis: key lemmas



min freq = 5, min range = 5
(both corpora)

NB how they differ
NOT what they have in common

min freq = 5, min range = 5
(both ways)

keylemmas 2018-2022
vs 1997-2017

keylemmas 1997-2017
vs 2018-2022

**1997-2017**

| | | | |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | that |
| 24 | book | 49 | reading |
| 25 | esl | 50 | micase |

**themes**

**going down the 1997-2017 list**
(principle uses)

**plus key n-grams**
(AntConc v3 workaround)

**2018-2022**

| | | | |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

# Corpus analysis: 50 key lemmas

## 1997-2017

| # | lemma | # | lemma |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | reading |
| 24 | book | 49 | micase |
| 25 | esl | 50 | cobuild |

## technology, tools

| 1997-2017 | 2018-2022 |
|---|---|
| concordancer | ddl |
| concordance | query, line |
| concordancing | |
| operation | |
| occurrence, gram | |
| web, computer, google resource, suite, checker | app, mobile, platform tool |
| bank, bnc, micase, cobuild, ldoce book, text | skell |

## 2018-2022

| # | lemma | # | lemma |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

# Corpus analysis: 50 key lemmas

## 1997-2017

| | | | |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | reading |
| 24 | book | 49 | micase |
| 25 | esl | 50 | cobuild |

## 2018-2022

| | | | |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

## people involved

| 1997-2017 | 2018-2022 |
|---|---|
| student, trainee | learner, phd |
| ns | teacher |

# Corpus analysis: 50 key lemmas

## 1997-2017

| # | | # | |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | reading |
| 24 | book | 49 | micase |
| 25 | esl | 50 | cobuild |

## 2018-2022

| # | | # | |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

## (language) focus

| 1997-2017 | 2018-2022 |
|---|---|
| legal, glossary, grammar particle, que, stance routine, sequence | pronunciation, fluency error, correction corrective, thesis complexity, variation vocabulary |
| translation, parallel interpreting writer, reading | ra |
| esl, french, german | mandarin, cantonese arabic, hong kong |

# Corpus analysis: 50 key lemmas

## 1997-2017

| | | | |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | reading |
| 24 | book | 49 | micase |
| 25 | esl | 50 | cobuild |

## 2018-2022

| | | | |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

## activities

| 1997-2017 | 2018-2022 |
|---|---|
| project, module | workshop, submission lesson, blend |
| problem, scaffolding | |
| conceptual, procedural focal | |
| exercise, gloss, cloze | instruction |

**1997-2017**

| | | | |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | reading |
| 24 | book | 49 | micase |
| 25 | esl | 50 | cobuild |

some things changed
some things disappeared…

what's completely new?

themes
continuing down the 2018-2022 list

**2018-2022**

| | | | |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

# Corpus analysis: 50 key lemmas

## 1997-2017

| | | | |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | reading |
| 24 | book | 49 | micase |
| 25 | esl | 50 | cobuild |

## 2018-2022

| | | | |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

## areas

| 1997-2017 | 2018-2022 |
|---|---|
| | retention, memory |
| | anxiety, enjoyment |

## research

| 1997-2017 | 2018-2022 |
|---|---|
| | al, et, boers, lee, webb |
| | crosthwaite |
| | iteration, rq, pre, post |
| | min, covariate, effect |

# Corpus analysis: key n-grams 2018-23

| Rank | 3-gram (x270) |
|------|---------------|
| 1 | lee et al |
| 2 | in hong kong |
| 3 | corpus based tasks |
| 4 | of the target |
| 5 | boulton and cobb |
| 6 | the effect of |
| 7 | and learner corpora |
| 8 | and post tests |
| 9 | the pre and |
| 10 | as shown in |
| 11 | pre and post |
| 12 | in terms of |
| 13 | of corpus tools |
| 14 | the effectiveness of |
| 15 | immediate and delayed |
| 16 | the number of |
| 17 | the post test |
| 18 | and genre based |
| | et al p |
| | of the error |
| | the target collocations |
| 22 | the pre test |
| 23 | the concordance lines |
| 24 | the present study |
| 25 | the participants of |

| Rank | 4-gram (x140) |
|------|---------------|
| 1 | pre and post tests |
| 2 | the pre and post |
| 3 | data driven learning ddl |
| 4 | in the pre and |
| 5 | the pre test to |
| 6 | use of corpus tools |
| 7 | in the post test |
| 8 | corpus of contemporary american |
| | of contemporary american english |
| 10 | as a learning tool |
| | between the pre and |
| | to use corpus tools |
| 13 | as shown in table |
| 14 | effects of ddl on |
| | students awareness of the |
| | the control and experimental |
| 17 | in english language teaching |
| 18 | as a foreign language |
| 19 | raise students awareness of |
| | the meanings of the |
| | the participants of the |
| | to be more effective |
| | in the pre test |
| 24 | in terms of the |
| 25 | engine for language learning [+10] |

| Rank | 5-gram (x35) |
|------|---------------|
| 1 | the pre and post tests |
| 2 | the use of corpus tools |
| 3 | corpus of contemporary american english |
| 4 | in the pre and post |
| 5 | between the pre and post |
| | the control and experimental groups |
| 7 | sketch engine for language learning |
| 8 | the corpus of contemporary american |
| 9 | findings of the present study |
| | the long term effects of |
| | the pre test and post |
| 12 | the use of the corpus |
| 13 | on the basis of the |
| 14 | in the pre test and |
| | it is worth mentioning that |
| | of english for academic purposes |
| | the findings of the present |
| | to be more effective than |
| 19 | english as a foreign language |
| 20 | english for international communication toeic |
| | of language learning and teaching |
| | participants were randomly divided into |
| | the present study aims to |
| | the test of english for |
| 25 | data driven learning ddl johns |

# Corpus analysis: 50 key lemmas

## 1997-2017

| # | lemma | # | lemma |
|---|---|---|---|
| 1 | concordancer | 26 | resource |
| 2 | student | 27 | operation |
| 3 | web | 28 | checker |
| 4 | ns | 29 | sequence |
| 5 | concordance | 30 | exercise |
| 6 | legal | 31 | french |
| 7 | grammar | 31 | bnc |
| 8 | concordancing | 31 | text |
| 9 | translation | 34 | conceptual |
| 10 | project | 35 | trainee |
| 11 | parallel | 36 | occurrence |
| 12 | stance | 37 | glossary |
| 13 | bank | 38 | german |
| 14 | computer | 39 | gloss |
| 15 | problem | 39 | scaffolding |
| 16 | ldoce | 41 | focal |
| 16 | particle | 42 | grasp |
| 18 | example | 43 | module |
| 19 | writer | 44 | procedural |
| 20 | que | 45 | gram |
| 21 | google | 46 | suite |
| 22 | interpreting | 47 | cloze |
| 23 | routine | 48 | reading |
| 24 | book | 49 | micase |
| 25 | esl | 50 | cobuild |

## 2018-2022

| # | lemma | # | lemma |
|---|---|---|---|
| 1 | pronunciation | 26 | rq |
| 2 | fluency | 27 | lesson |
| 3 | error | 28 | lee |
| 4 | app | 29 | vocabulary |
| 5 | workshop | 30 | arabic |
| 6 | skell | 31 | post |
| 7 | teacher | 31 | min |
| 8 | al | 31 | cantonese |
| 9 | learner | 34 | webb |
| 10 | anxiety | 35 | foreign |
| 11 | variation | 36 | hong |
| 12 | correction | 37 | kong |
| 13 | et | 38 | covariate |
| 14 | mobile | 39 | ra |
| 15 | ddl | 39 | pre |
| 16 | instruction | 41 | line |
| 17 | phd | 42 | learning |
| 18 | enjoyment | 43 | query |
| 19 | retention | 44 | crosthwaite |
| 20 | platform | 45 | complexity |
| 21 | submission | 46 | corrective |
| 22 | mandarin | 47 | blend |
| 23 | thesis | 48 | effect |
| 24 | iteration | 49 | memory |
| 25 | boers | 50 | tool |

### concordancing

| 1997-2017 | | 2018-2022 |
|---|---|---|
| 757 | frequency | 120 |
| 0.97 | per thousand words | 0.33 |
| 83/122 (68.0%) | range | 21/53 (39.6%) |

### DDL

| 1997-2017 | | 2018-2022 |
|---|---|---|
| 1874 | frequency | 1345 |
| 2.41 | per thousand words | 3.74 |
| 56/122 (45.9%) | range | 42/53 (79.2%) |

# Other corpus tools (DDL)



1997-2003 | 2006-2011 | 2011-2014 | 2014-2016 | 2016-2018 | 2018-2020 | 2020-2022

## DDL over time (inc 0 mentions)

### Top DDL RAs ptw

| | Date | RAFreq | ptw |
|---|---|---|---|
| 1 | 2016 | Mizumoto & Chujo | 29.11 |
| 2 | 2015 | Lin & Lee | 27.48 |
| 3 | 2016 | Mizumoto et al | 22.46 |
| 4 | 2019 | Lin & Lee | 21.45 |
| 5 | 2016 | Lin | 19.94 |
| 6 | 2020 | Saeedakhtar et al | 17.44 |
| 7 | 2020 | Lee et al | 16.13 |
| 8 | 2016 | Vyatkina (a) | 15.61 |
| 9 | 2016 | Vyatkina (b) | 15.35 |
| 10 | 2018 | Moon & Oh | 14.72 |
| 11 | 2014 | Smart | 14.62 |
| 12 | 2016 | Karras | 14.61 |
| 13 | 2021 | Gilquin | 13.58 |
| 14 | 2010 | Boulton | 13.41 |
| 15 | 2017 | Ackerley | 12.35 |
| 16 | 2022 | Samoudi & Modir. | 11.50 |
| 17 | 2017 | Hadley & Charles | 10.61 |
| 18 | 2019 | Crosthwaite et al | 10.06 |

## DDL: ptw



## DDL: context

(cf. Boulton, 2011)

Johns 1986. *Micro-Concord: A language learner's research tool.*
- "concordancing"

Johns 1988. *Whence and whither classroom concordancing?*

Johns 1990. *From printout to handout: Grammar and vocabulary teaching in the context of data-driven learning.*

Johns 1991. *Should you be persuaded: Two samples of data-driven learning.*

Johns & King (eds.) 1991. *Classroom Concordancing.*
- "an application of computers to language-learning that has come to be known as 'classroom concordancing' or 'data-driven learning' (DDL)" (p.iii)

Johns 1993. *Data-driven learning: An update.*
- "The earlier term Classroom Concordancing described the technique; the new term Data-Driven Learning was coined to emphasise the methodology." (p.4)

Johns 2002. *Data-driven learning: The perpetual challenge.*
- "an approach… that I have, for want of a better term, named data-driven learning." (p.107)

Johns et al. 2008. *Integrating corpus-based CALL programs in teaching English through children's literature.*
- "corpus-based language learning" (p.495)

175 comparison texts: ±ISLA ⇨ same journal, same year (same issues); min range = 5

| Rank | DDL keywords (x759) | | |
|---|---|---|---|
| 1 | corpus | 26 | chambers |
| 2 | corpora | 27 | concordancers |
| 3 | ddl | 28 | patterns |
| 4 | concordance | 29 | deductive |
| 5 | collocations | 30 | bnc |
| 6 | collocation | 31 | word |
| 7 | concordancing | 32 | driven |
| 8 | concordancer | 33 | google |
| 9 | search | 34 | materials |
| 10 | concordances | 35 | data |
| 11 | consultation | 36 | linguistics |
| 12 | boulton | 37 | formulaic |
| 13 | use | 38 | based |
| 14 | searches | 39 | coca |
| 15 | query | 40 | noun |
| 16 | examples | 41 | lexico |
| 17 | johns | 42 | errors |
| 18 | lines | 43 | reference |
| 19 | cobb | 44 | phrases |
| 20 | tools | 45 | collocational |
| 21 | inductive | 46 | approach |
| 22 | hands | 47 | kennedy |
| 23 | yoon | 48 | verb |
| 24 | collocates | 49 | exercises |
| 25 | queries | 50 | preposition |

| Rank | Non-DDL keywords (x1048) | | |
|---|---|---|---|
| 1 | captions | 26 | planning |
| 2 | interaction | 27 | spanish |
| 3 | social | 28 | chat |
| 4 | communication | 29 | multimodal |
| 5 | technology | 30 | vowel |
| 6 | face | 31 | cmc |
| 7 | feedback | 32 | wiki |
| 8 | captioning | 33 | cultural |
| 9 | collaborative | 34 | video |
| 10 | self | 35 | exchange |
| 11 | comprehension | 36 | cf |
| 12 | peer | 37 | digital |
| 13 | negotiation | 38 | voice |
| 14 | l | 39 | graph |
| 15 | listening | 40 | blended |
| 16 | mall | 41 | call |
| 17 | messages | 42 | strategies |
| 18 | environment | 43 | game |
| 19 | mail | 44 | facebook |
| 20 | mobile | 45 | emotions |
| 21 | scmc | | synchronous |
| 22 | practices | 47 | global |
| 23 | clil | 48 | virtual |
| | blog | 49 | reading |
| 25 | captions | 50 | semiotic |

Different syntheses (NS, MA, MM, corpus): complementary, triangulation
  ⇨ essential to know your field! Automated, statistics, but…

1. Listen to past recommendations: better research practices,
   greater rigor in reporting (e.g. duration, proficiency, activities, materials)

2. More diversity, originality
   'corpus' types, tools & interfaces
   AI/ChatGPT?

3. Research on the underpinnings of DDL (processes), e.g. DDL promotes
   autonomy, noticing, induction, language awareness …'better learners'?

http://micase.elicorpora.info/

Analyse1000100101000100110001101010100011
010011et010100110001110011101010010101011
1Traitement01010001100010101011001101
01001Informatique0101001011001001
de0101la0100011101010001
0101Langue010111001
Francaise01010010
Conversation
Analyse

your program there's a form that you can mail in. um **thank you** and have a wonderful evening. APPLAUSE {END

s i learned to analyze scientific research articles later. **thank you** for not making me dread that. you have the ability to

**thank you** have a nice weekend UNINTELLIGIBLE CONVERSATION

uh so let's give the tape recorder a break too, and so **thank you** very much and i'll see you on Thursday. micase-related

PAUSE duration well thank you. **thank you**. {END OF TRANSCRIPT}

of getting the slides please? okay. uh there we go um, **thank you**. now look at his, another image of augustus here, um this

down and if you could just give them over to nikolas. **thank you**. UNINTELLIGIBLE SPEECH

uh i think we can just actually, stop slides, yeah **thanks** (we can get) a little more light here. um, and the scale of this

ion, right through those doors to the right, afterwards, **thanks**, for coming everyone APPLAUSE {END OF TRANSCRIP

ions before we wrap up...? okay, that concludes it then **thanks**. {END OF TRANSCRIPT}

ery nice. any questions? okay **thanks**. okay. all righty um what I

, highly intensive coffee plantations. SLIDE CHANGE **thanks**. so, given this context, then uh obviously one of the things

's a bunch of extras here. oh **thanks**. PAUSE WHILE LOWERING SCREEN so, again this is one of

m, well thank you very much. i think we're done, and **thanks** for, allowing this to be videotaped, this project thanks you.

wrry about things that we haven't discussed at all. so, **any questions** (coming up?) everyone's is th

cover on aquifer evaluation tests but i i are there **any questions**? is everyone i, you can't learn all the all the details

ore we get going with the selection sort again are there **any questions** about anything...? okay. well what i'd like to do first...

o do the exchange in the other array. kay well are there **any questions** about this? PAUSE duration :05 kay well let's start

i'm gonna assign uh practice problems for homework. **any questions** before we wrap up...? okay, that concludes it then

stions you, make sure that if you have any concerns, **any questions** email me. and what would be better is if you can

yes i'll entertain **any questions** i'm dying to ask you a questi