



**Recueil, traitement et analyse des corpus d'apprenants : Défis méthodologiques et enjeux**

Gaëtanelle Gilquin,  
Université catholique de Louvain

JECA, 18 mars 2022

1

Quelle est la première question à se poser avant de commencer à constituer un corpus (d'apprenants) ?

“Existe-il un corpus qui répond à mes besoins ?”

2

**Une multitude de corpus d'apprenants**

- Les corpus d'apprenants ont commencé à apparaître au début des années 1990 (*International Corpus of Learner English, Longman Learners' Corpus*)
- Depuis lors, de nombreux corpus d'apprenants ont été constitués, dont certains ont été rendus accessibles à la communauté
- Majorité de corpus écrits et de corpus d'anglais (L2), mais d'autres types de corpus disponibles

3

Learner corpora <https://uclouvain.be/en/research-institutes/ice/cecl/learner-corpora-around-the-world.html>

Use the query box below to search for specific keywords (e.g. languages, task-type, medium).

Show: All entries Search: French written

Corpus	Target language	First language	Medium	Text type / task type	Proficiency level	Size in words	Project director	Availability
The "Dile Automatism" corpus	French	Mainly L1 speakers of English	Written	Narrative, persuasive, persuasive and informative texts		c. 50,000	Marie Josée Hamel, Suzanne Mottier (Dalhousie University, Canada)	Available after registration
French Interlanguage Database (FRIDAN)	French	Various	Written	Free compositions: descriptive, argumentative and narrative texts, news & mail; Argumentative essays, informative texts, journalistic texts, formal letters, summaries, written compositions by French students of French	Intermediate		Sylvain Granger (Centre for English Corpus Linguistics, Université catholique de Louvain, Belgium)	
The Learner Corpus French (LCF)	French	Dutch	Written	Argumentative essays, summaries, written compositions by French students of French	Intermediate to advanced	c. 500,000	Hans Poelsman (KU Leuven/UGent / Leuven, Belgium)	Under development
The Lund CEFL-E Corpus (Corpus Ecrit de Français Langue Étrangère)	French	Swedish	Written	Descriptive and narrative essays, lecture-based dialogues	Various	c. 100,000	Matti Agren (Lund University, Sweden)	A sub-part of the corpus is available online.

4

**La constitution de corpus d'apprenants**

- Nécessaire si le corpus souhaité n'existe pas encore, par exemple :
  - langue (L1/L2) pas/peu représentée
  - variété très spécifique (ex. production écrite d'apprenants de l'anglais en pharmacologie)
  - variété locale (ex. travaux de ses propres étudiants de Master)
- Importance de la constitution des corpus : 'Garbage in, Garbage out' (cf. Sinclair 1991, p.9)

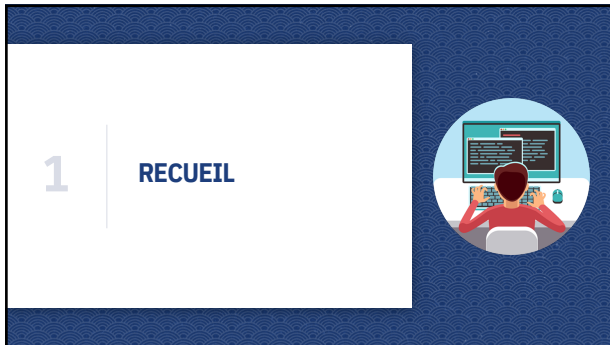
5

**Les étapes**

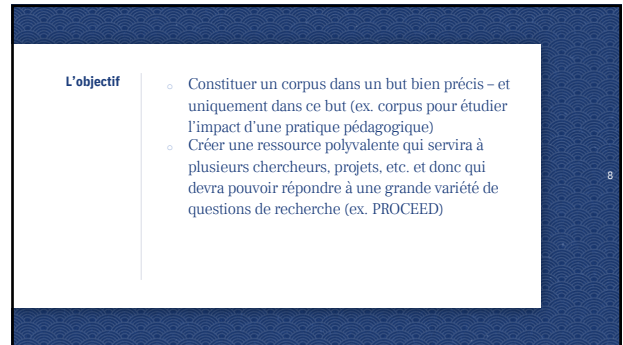
1. Recueil
2. Traitement
3. Analyse

- Chacune de ces étapes peut représenter des défis méthodologiques majeurs (accent sur ce qui fait la spécificité des corpus d'apprenants)
- Illustration avec certains des projets de corpus d'apprenants initiés au CECL : ICLE, ICLE+30, LONGDALE, PROCEED, LINDSEI
- Surtout l'anglais, mais applicable à d'autres langues

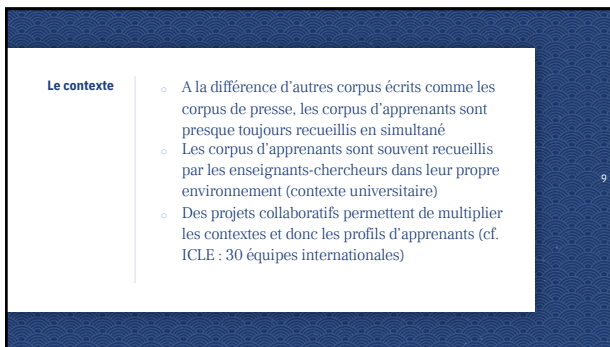
6



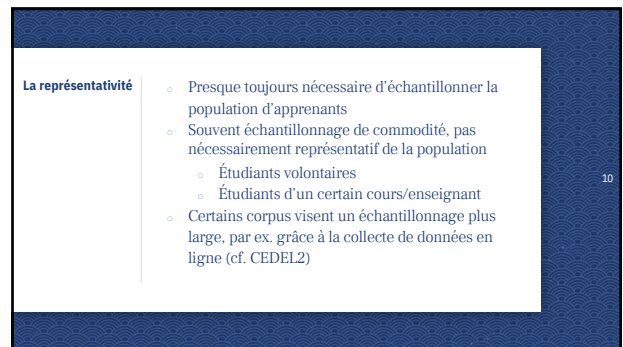
7



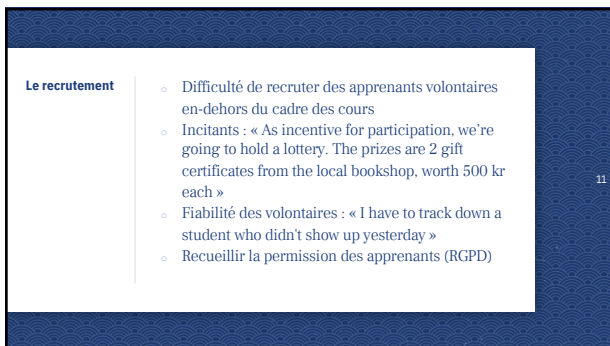
8



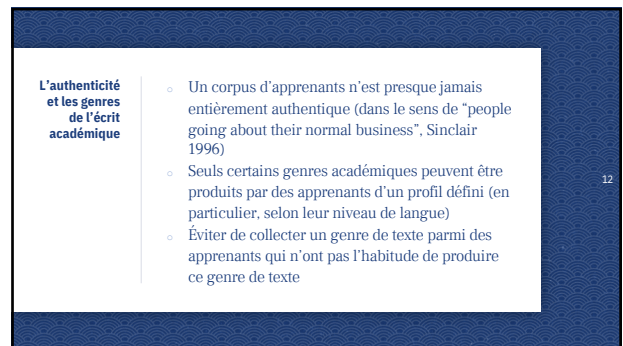
9



10



11



12

**Le cas des corpus longitudinaux**

- Projets de longue durée (financement, personnel)
- Perte de participants (attrition)

Trajectory	Number of students
Y1	86
Y1-Y2	17
Y1-Y2-Y3	66
Y1-Y3	15
Y2	10
Y2-Y3	16
Y3	27
<b>Total</b>	<b>237</b>

LONGDALE (Paquot et al. 2021, p.129)

13

**Les métadonnées**


- Indispensable de collecter des métadonnées riches et détaillées, étant donné l'hétérogénéité des profils d'apprenants et l'influence de nombreuses variables sur la langue d'apprenants (Ellis 1994, Granger 1998)
- Sélection de variables potentiellement pertinentes, cf. Lozano (2021) (ex. interviewer dans LINDSEI)
- Niveau de compétence : si possible inclure le résultat d'un test standardisé, ex. LexTALE (limite des indicateurs de substitution, cf. Myles 2015, p.316)
- Impossible d'inclure toutes les variables pouvant influencer la langue d'apprenants (Granger 2004, p.126)

14

“Corpus-building is of necessity a marriage of perfection and pragmatism” (McEnery et al. 2006, p.73)

15

2 **TRAITEMENT**

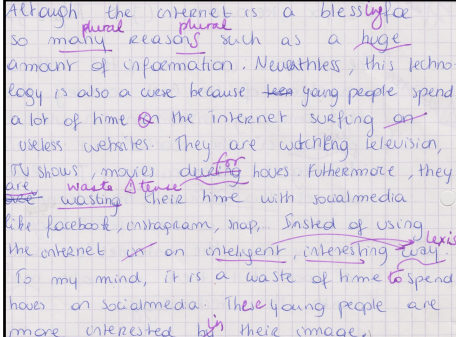


16

**La numérisation**

- Aujourd'hui, de nombreux corpus d'apprenants rassemblent des textes produits directement sous format électronique (cf. ICLE+30)
- Certains contextes où l'écriture à la main est encore de mise (ex. apprenants débutants, examens) => travail manuel de transposition
- Crucial de reproduire le texte exactement à l'identique - y compris les formes non-standard

17



Although the internet is a blessing for so many <sup>plural</sup> reasons <sup>plural</sup> such as a huge amount of information. Nevertheless, this technology is also a curse because ~~teen~~ young people spend a lot of time on the internet surfing on useless websites. They are watching television, TV shows, movies ~~during~~ hours. Furthermore, they are ~~waste~~ <sup>wasting</sup> their time with socialmedia like facebook, instagram, snapchat. Instead of using the internet as an intelligent, interesting <sup>lexis</sup> way. To my mind, it is a waste of time to spend hours on socialmedia. These young people are more interested <sup>by</sup> their image.

18





**L'annotation (2)**

- Résultats de l'étiquetage morpho-syntaxique et de l'analyse syntaxique de corpus d'apprenants généralement assez satisfaisants
  - Structure relativement simple des phrases (Meunier 1998, Huang et al. 2018)
  - Meilleurs résultats pour des niveaux de compétence avancés (Geertzen et al. 2014)
  - Peut varier selon les textes : 95%-99.1% de précision avec CLAWS sur ICLE (Granger et al. 2009, p.16)
  - Peut varier selon les catégories, ex. 100% pour JJR vs 33% pour VBO avec CLAWS sur LINDSEI

25

**L'annotation (3): les erreurs**

- L'annotation automatique d'erreurs est pour l'instant limitée à des phénomènes spécifiques, ex. prépositions (De Felice & Pulman 2009), articles (Rozovskaya & Roth 2010), orthographe (Rayson & Baron 2011)
- Comment identifier les erreurs ?
  - Études pilotes autour de l'identification d'erreurs par des locuteurs natifs : parfois plus du double d'erreurs identifiées par texte (Dagneaux et al. 2008, p.7)

26

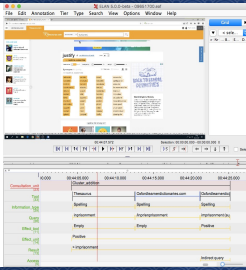
**L'annotation (4): les erreurs**

- Comment distinguer les erreurs de "maladresses" ?
  - Louvain Error Tagging Manual (Dagneaux et al. 2008) : 'Errors' vs 'Infelicities' (ex. couldn't \$could not\$ ; him \$him/her\$) – sauf connecteurs et ponctuation ; quid des collocations (ex. high \$heavy\$ responsibilities) ?
- Comment déterminer la meilleure correction ?
  - Sometimes we cannot blame the masses for being (LS) narrow \$small-minded?\$...
  - If we were to assume that technological progress and human imagination are two irreconcilable things, then we would (LS) imminently \$?\$ be taking on a Romantic or even Luddite attitude.

27

**L'annotation (5): autres**

- Toute couche d'annotation augmente la valeur du corpus et permet de répondre à davantage de questions
- Impossible d'inclure tous les types d'annotation potentiellement utiles => à ajouter selon les besoins
- Process Corpus of English in Education (PROCEED : Gilquin 2022a) : annotation de l'utilisation d'outils d'aide à la rédaction (Gilquin & Laporte 2021)




28

“ In practice, most annotation schemes ... aim at the maximum potential utility, tempered by the practicalities of annotating the text” (McEnery & Wilson 2001, p.34)

29

3 ANALYSE



30

**L'objet de recherche (1): corpus-driven**

- Privilégier une approche centrée sur le corpus / la langue d'apprenants, ex.
- Recherche du mot *library* : 40 orthographes différentes dans un corpus d'anglais d'apprenants japonais (Milton & Okada 2007) : *libelary, library, liburality, liburary, liveraly, liverary, liverely, etc.* (=>VariAnt, Anthony 2017)
- Recherche de connecteurs sur base d'une liste pré-établie (Granger & Tyson 1996) => quid des connecteurs mal orthographiés (ex. *although, inspite, therefor*) ou non-standard (ex. *on the other side, according to me*) ?
- "Recall problem" / "You-Don't-Know-What-You're-Missing problem" (Ball 1994, p.295)

31

**L'objet de recherche (2): annotation**

- Commencer par une étude pilote pour évaluer la précision de la recherche automatique (Granger 1997)
- Recherche du passif (Granger 1997), ex.
  - \* *The reader is lead-noun> to believe that...*
  - \* *The laws aren't uphold-Vin/> that efficiently*
- Recherche des constructions causatives avec *cause, get, have et make* (Gilquin 2016)
  - \* Inclusion de toutes les étiquettes correspondant à une forme verbale non-personnelle (infinitif, participe passé et présent), ex. *cause to decrease the crime rate ; makes people trying to find one*

32

**Les facteurs d'influence (1)**

- Les nombreux éléments du profil socio-linguistique de l'apprenant (et leur interaction !)
- Souvent impossible d'effectuer un contrôle strict dans la sélection des données à cause de la taille du corpus (Callies 2015, p.52)
- Interactions parfois difficiles à interpréter, cf. Deshors (2021) à propos des verbes à particule en anglais dans VOICE : type d'objet direct X rôle de l'orateur

33

**Les facteurs d'influence (2)**

- La tâche, y compris l'intitulé / le sujet

Paquot (2013, p.399) : 14 des 34 blocs lexicaux distinctifs de ICLE-FR (parmi ICLE) sont liés au sujet de la création de l'Europe, ex.

- *Europe will be united against USA and Japan.*
- *Each country will keep its own identity, currency, institutions and constitution.*

Bell et al. (2021, p.223-4) : Image de deux policiers, d'un enfant et de sa mère

- Moins d'opportunités d'utiliser des mots féminins que masculins
- Incertitude quant au caractère correct ou incorrect des pronoms décrivant l'enfant

34

**Les facteurs d'influence (3)**

- La longueur des textes
  - La longueur des phrases tend à augmenter proportionnellement à la longueur des textes (Buttery & Caines 2012)
  - Les adverbes ont proportionnellement plus de chances d'être utilisés dans des textes plus longs (ibid.)
  - La nature et la fréquence des blocs lexicaux varie selon la longueur des textes (et le nombre de textes par corpus) (Pan et al. 2020)

35

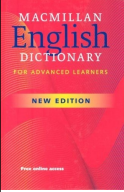
**Le choix de la norme (1) (Gilquin 2022b)**

- Étudiants natifs ?
  - "Native-speaking students do not necessarily provide models that everyone would want to imitate" (Leech 1998, p.xix)
- Experts (natifs ou non-natifs) ?
  - "both unfair and descriptively inadequate" (Lorenz 1999, p.14)
  - "unrealistic standard of 'expert writer' models" (Hyland & Milton 1997, p.184)
- Étudiants natifs si évaluation ; experts si but formatif (Ádel 2006, p.206-207)

36

**Le choix de la norme (2)**

- o Ressource pédagogique pour aider les apprenants à améliorer leurs compétences rédactionnelles en anglais académique
  - o 100 notes 'Get it right' sous les entrées du dictionnaire correspondantes
  - o *Improve your writing skills*: 12 sections axées sur des fonctions rhétoriques importantes en anglais académique (Gilquin et al. 2007a, 2007b)
  - o Basée sur l'analyse d'ICLE et sa comparaison avec LOCNESS, mais surtout avec le BNC



37

**Improve your writing skills**

**1. Giving your opinion explicitly**

**2. Giving your opinion implicitly**

**3. Expressing possibility and tentativeness**

**4. Using modal auxiliaries with verbs**

**5. Using the subjunctive**

**6. Using the imperative**

**7. Using the infinitive**

**8. Using the gerund**

**9. Using the participle**

**10. Using the relative clause**

**11. Using the conditional**

**12. Using the passive voice**

**13. Using the passive voice in the infinitive**

**14. Using the passive voice in the gerund**

**15. Using the passive voice in the participle**

**16. Using the passive voice in the relative clause**

**17. Using the passive voice in the conditional**

**18. Using the passive voice in the infinitive**

**19. Using the passive voice in the gerund**

**20. Using the passive voice in the participle**

**21. Using the passive voice in the relative clause**

**22. Using the passive voice in the conditional**

**23. Using the passive voice in the infinitive**

**24. Using the passive voice in the gerund**

**25. Using the passive voice in the participle**

**26. Using the passive voice in the relative clause**

**27. Using the passive voice in the conditional**

**28. Using the passive voice in the infinitive**

**29. Using the passive voice in the gerund**

**30. Using the passive voice in the participle**

**31. Using the passive voice in the relative clause**

**32. Using the passive voice in the conditional**

**33. Using the passive voice in the infinitive**

**34. Using the passive voice in the gerund**

**35. Using the passive voice in the participle**

**36. Using the passive voice in the relative clause**

**37. Using the passive voice in the conditional**

**38. Using the passive voice in the infinitive**

**39. Using the passive voice in the gerund**

**40. Using the passive voice in the participle**

**41. Using the passive voice in the relative clause**

**42. Using the passive voice in the conditional**

**43. Using the passive voice in the infinitive**

**44. Using the passive voice in the gerund**

**45. Using the passive voice in the participle**

**46. Using the passive voice in the relative clause**

**47. Using the passive voice in the conditional**

**48. Using the passive voice in the infinitive**

**49. Using the passive voice in the gerund**

**50. Using the passive voice in the participle**

**51. Using the passive voice in the relative clause**

**52. Using the passive voice in the conditional**

**53. Using the passive voice in the infinitive**

**54. Using the passive voice in the gerund**

**55. Using the passive voice in the participle**

**56. Using the passive voice in the relative clause**

**57. Using the passive voice in the conditional**

**58. Using the passive voice in the infinitive**

**59. Using the passive voice in the gerund**

**60. Using the passive voice in the participle**

**61. Using the passive voice in the relative clause**

**62. Using the passive voice in the conditional**

**63. Using the passive voice in the infinitive**

**64. Using the passive voice in the gerund**

**65. Using the passive voice in the participle**

**66. Using the passive voice in the relative clause**

**67. Using the passive voice in the conditional**

**68. Using the passive voice in the infinitive**

**69. Using the passive voice in the gerund**

**70. Using the passive voice in the participle**

**71. Using the passive voice in the relative clause**

**72. Using the passive voice in the conditional**

**73. Using the passive voice in the infinitive**

**74. Using the passive voice in the gerund**

**75. Using the passive voice in the participle**

**76. Using the passive voice in the relative clause**

**77. Using the passive voice in the conditional**

**78. Using the passive voice in the infinitive**

**79. Using the passive voice in the gerund**

**80. Using the passive voice in the participle**

**81. Using the passive voice in the relative clause**

**82. Using the passive voice in the conditional**

**83. Using the passive voice in the infinitive**

**84. Using the passive voice in the gerund**

**85. Using the passive voice in the participle**

**86. Using the passive voice in the relative clause**

**87. Using the passive voice in the conditional**

**88. Using the passive voice in the infinitive**

**89. Using the passive voice in the gerund**

**90. Using the passive voice in the participle**

**91. Using the passive voice in the relative clause**

**92. Using the passive voice in the conditional**

**93. Using the passive voice in the infinitive**

**94. Using the passive voice in the gerund**

**95. Using the passive voice in the participle**

**96. Using the passive voice in the relative clause**

**97. Using the passive voice in the conditional**

**98. Using the passive voice in the infinitive**

**99. Using the passive voice in the gerund**

**100. Using the passive voice in the participle**

38

**Le circuit long**

- o Corpus pour utilisation pédagogique différée (Granger 2004)
- o Risque que le circuit soit interrompu car acteurs différents

39

**Le circuit court**

- o Corpus pour utilisation pédagogique immédiate (Granger 2004)
- o 'Corpus d'apprenants local' (Seidhofer 2002)
- o Ampleur limitée, mais sur mesure

40

**L'évolution diachronique**

- o Date de péremption des corpus d'apprenants ?
  - o Contextes d'acquisition
  - o Mode d'écriture
  - o Format des corpus
- o ICLE+30 (cf. Gilquin 2021, à paraître)
  - o Recueil de données similaires à ICLE 30 ans plus tard
  - o Données plus actuelles
  - o Analyse diachronique de la langue d'apprenants en comparant ICLE et ICLE+30

41

“Every corpus I have had the chance to examine, however small, has taught me facts I couldn't imagine finding out any other way” (Fillmore 1992, p.35)

42

# 4 CONCLUSION

43

- Recueillir et traiter son corpus pour mieux l'analyser**
- Les défis et les difficultés liés au recueil, au traitement et à l'analyse des corpus d'apprenants sont nombreux, mais...
  - La constitution de corpus d'apprenants est une expérience enrichissante
  - Constituer et traiter un corpus permet de mieux connaître ses données, ce qui peut faciliter leur exploitation
  - Service à la communauté, mutualisation des ressources

44

- Quelques lectures utiles**
- Bell, Philippa & Caroline Payant. 2021. Designing learner corpora: Collection, transcription, and annotation. In Nicole Tracy-Ventura & Magali Paquot (eds) *The Routledge Handbook of Second Language Acquisition and Corpora* (pp. 53-67). Abingdon: Routledge.
  - Gilquin, Gaëtan. 2015. From design to collection of learner corpora. In Sylviane Granger, Gaëtan Gilquin & Fanny Meunier (eds) *The Cambridge Handbook of Learner Corpus Research* (pp. 9-34). Cambridge: Cambridge University Press.
  - Tono, Yukio. 2016. What is missing in learner corpus design? In Margarita Alonso-Ramos (ed.) *Spanish Learner Corpus Research: Current Trends and Future Perspectives* (pp. 33-52). Amsterdam: John Benjamins.

45

**Références (1)**

- Ald, A. 2006. *Metascourse in L1 and L2 English*. Amsterdam, John Benjamins.
- Anthony, L. 2017. *VarXact* (Version 1.1.0) [Computer Software]. Tokyo, Japan: Waseda University. Available from <https://www.laurenceanthony.net/software>
- Ball, C. N. 1984. Automated text analysis: Cautionary tales. *Literary and Linguistic Computing* 9(4): 295-302.
- Belli, P., L. Collins & E. Manders. 2021. Building an oral and written learner corpus of a school programme: Methodological issues. In B. Le Bruyn & M. Paquot (eds) *Learner Corpus Research: More Second Language Acquisition* (pp. 214-242). Cambridge: Cambridge University Press.
- Bunney, P. & A. Cairns. 2012. Normalising frequency counts to account for opportunity of use in learner corpora. In Y. Tono, Y. Kawaguchi & M. Minegishi (eds) *Developmental and Crosslinguistic Perspectives on Learner Corpus Research* (pp. 187-204). Amsterdam: John Benjamins.
- Callis, M. 2015. Learner corpus methodology. In S. Granger, G. Gilquin & F. Meunier (eds) *The Cambridge Handbook of Learner Corpus Research* (pp. 35-56). Cambridge: Cambridge University Press.
- Chapman, E., S. Desnos, S. Granger, F. Meunier, J. Nef & J. Thewissen. 2008. *The Learner error tagging manual: Version 1.3*. Centre for English Corpus Linguistics, Université catholique de Louvain.
- de Felice, R. & S. Palmu. 2009. Automatic detection of preposition errors in learner writing. *CALLA Journal* 26(3): 512-528.
- Dehors, S. 2021. Plural verbs in English as a Lingua Franca: a case for corpus-based multifactorial analysis. Presentation at *Workshop on Multivert Units in Multilingual Spaces*, Eberhard Karls Universität Tübingen, Germany, 7 June 2021.
- Ellis, B. 1994. *The Study of Second Language Acquisition*. Oxford: Oxford University Press.
- Geertman, J., T. Alexopoulos & A. Kerfoot. 2014. Automatic linguistic annotation of large scale L2 databases: The IP-Cambridge Open Language Database (IP-CamLab). In B. T. Miller & M. C. M. Edgington, A. Henry, N. Marcos Miguel, A. M. Terey, A. Tuzimetti, & H. Walter (eds) *Solved Proceedings of the 2012 Second Language Research Forum: Building Bridges between Disciplines* (pp. 240-254). Somerville: Cascadia Proceedings Project.
- Hillemann, C. J. 1992. "Corpus linguistics" or "Computational grammar linguistics". In J. Swainik (ed.) *Directions in Corpus Linguistics. Proceedings of Applied Symposium 52, Stockholm, 4-8 August 1991* (pp. 35-60). Berlin: Mouton de Gruyter.
- Gilquin, G. 2016. Input-dependent L2 acquisition: Causative constructions in English as a foreign and second language. In S. De Knip & G. Gilquin (eds) *Applied Computational Grammar* (pp. 115-148). Berlin: de Gruyter.
- Gilquin, G. 2021. *Her-stud directions: Exploring some terra incognita in learner corpus research*. In A. Cermakova & M. Mali (eds) *Variation in Time and Space: Enriching the World through Corpora* (pp. 62-86). Berlin: De Gruyter.

46

**Références (2)**

- Gilquin, G. 2022a. The Process Corpus of English in Education: Going beyond the written text. *Research in Corpus Linguistics* 10(1): 31-44. Available at <http://red.scripps.edu/~gilqin/106-110-2022ILR1.pdf>
- Gilquin, G. 2022b. One norm to rule them all? Corpus-derived norms in learner corpus research and foreign language teaching. *Language Teaching* 55(1): 87-99.
- Gilquin, G. A parallel, diachronic learner corpus research: Examining learner language through the lens of time. In S. Flach & M. Hillert (eds) *Rethinking the Spectrum of Corpus Linguistics: New Approaches to Usability and Change*. Amsterdam: John Benjamins.
- Gilquin, G. & S. De Cock. 2011. Errors and disfluencies in spoken corpora: Setting the scene. *International Journal of Corpus Linguistics* 14(2): 141-172.
- Gilquin, G., S. Granger & M. Paquot. 2007a. Improve your writing skills (Writing sections). In M. Rundell (Editor in chief) *Macmillan English Dictionary for Advanced Learners (Second Edition)* (pp. 1011-1029). Oxford: Macmillan Education.
- Gilquin, G., S. Granger & M. Paquot. 2007b. Learner corpora: The missing link in EAP pedagogy. *Journal of English for Academic Purposes* 6(4): 319-335.
- Gilquin, G. & S. Laporte. 2021. The use of online writing tools by learners of English. *Lookout from a process corpus*. *International Journal of Corpus Linguistics* 14(4): 472-492.
- Granger, S. 1997. Automated retrieval of passives from native and learner corpora: precision and recall. *Journal of English Linguistics* 25(4): 365-374.
- Granger, S. 1998. The computer learner corpus: A versatile new source of data for SLA research. In S. Granger (ed.) *Learner English on Computer* (pp. 3-18). London: Longman.
- Granger, S. 2004. Computer learner corpus research: Current status and future prospects. In U. Connor & Th. A. Upton (eds) *Applied Corpus Linguistics: A Multidimensional Perspective* (pp. 123-143). Amsterdam: Rodopi.
- Granger, S., E. Desnos, F. Meunier & M. Paquot. 2009. *International Corpus of Learner English, Version 2: CD-ROM and Handbook*. Louvain-la-Neuve: Presses universitaires de Louvain.
- Granger, S. & S. Tyson. 1996. Connector usage in the English essay writing of native and non-native EFL speakers of English. *World Englishes* 15(1): 17-27.
- Huang, Y., A. Marikam, T. Alexopoulos & A. Korhonen. 2018. Dependency parsing of learner English. *International Journal of Corpus Linguistics* 21(1): 28-54.
- Jenkins, K. & J. Milton. 1997. Qualification and certainty in L1 and L2 students' writing. *Journal of Second Language Writing* 6(2): 183-216.
- Kooye, R. 2015. *The Marrying Corpus of International Learner English (MILE)*. In M. Callis & S. Gray (eds) *Learner Corpora in Language Testing and Assessment* (pp. 13-24). Amsterdam: John Benjamins.
- Loech, G. 1998. *Porfoco*. In S. Granger (ed.) *Learner English on Computer* (pp. 210-23). London: Longman.
- Lorenzo, F. 1999. *Adjective Intensionality - Learners versus Native Speakers: A Corpus Study of Argumentative Writing*. Amsterdam: Rodopi.

47

**Références (3)**

- Lozano, C. 2021. CELE2: Design, compilation and web interface of an online corpus for L2 Spanish acquisition research. *Second Language @ Research: Online First*. <https://doi.org/10.1177/0898264321101522>
- McEnery, T. & A. Wilson. 2001. *Corpus Linguistics: An Introduction, Second Edition*. Edinburgh: Edinburgh University Press.
- McEnery, T., B. Xiao & Y. Tono. 2006. *Corpus-Based Language Studies: An Advanced Resource Book*. London: Routledge.
- Meunier, F. 1998. Computer tools for interlanguage analysis: A critical approach. In S. Granger (ed.) *Learner English on Computer* (pp. 19-37). London: Longman.
- Mitsun, R. & T. Okada. 2007. The adaptation of an English spellchecker for Japanese writers. Paper presented at the *Symposium on Second Language Writing, 15-17 September 2007*, Nagoya, Japan. Available at <http://option.sok.ac.uk/9923/992.pdf>
- Myles, F. 2015. Second language acquisition theory and learner corpus research. In S. Granger, G. Gilquin & F. Meunier (eds) *The Cambridge Handbook of Learner Corpus Research* (pp. 308-322). Cambridge: Cambridge University Press.
- Pan, F., R. Bopper & D. Biber. 2020. Methodological issues in restrictive lexical bundle research: The influence of corpus design on bundle identification. *International Journal of Corpus Linguistics* 25(2): 215-229.
- Papay, M. 2012. Lexical bundles and L1 transfer effects. *International Journal of Corpus Linguistics* 15(3): 391-417.
- Papay, M., H. Nawa & S. Gray. 2021. Using syntactic co-occurrences to trace phonological complexity development in learner writing: Verb object structures in LINGMILE. In B. Le Bruyn & M. Paquot (eds) *Learner Corpus Research: More Second Language Acquisition* (pp. 122-147). Cambridge: Cambridge University Press.
- Rayson, P. & A. Barts. 2011. Automatic error tagging of spelling mistakes in learner corpora. In F. Meunier, S. De Cock, G. Gilquin, & M. Paquot (eds) *Taste for Corpora: An Honour to Sylviane Granger* (pp. 199-226). Amsterdam: John Benjamins.
- Roche, L., J. M. Hendricks & V. A. de Marjony. 2019. Estimating the success of re-identifications in incomplete datasets using generative models. *Nature Communications* 10: 3069.
- Shenoykova, A. & B. D. 2010. Training paradigms for correcting errors in grammar and usage. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics* (pp. 134-162). Los Angeles: Association for Computational Linguistics.
- Schiffelers, R. 2002. Pedagogy and local learner corpora: Working with learning-driven data. In S. Granger, J. Huang, & S. Peck-Tyson (eds) *Computer Learner Corpora, Second Language Acquisition and Foreign Language Teaching* (pp. 215-234). Amsterdam: John Benjamins.
- Schulz, J. 1991. *Corpus Construction: Calibration*. Oxford: Oxford University Press.
- Schulz, J. 1996. Preliminary recommendations on corpus pedagogy. Technical report, EAGLES (Expert Advisory Group on Language Engineering Standards). Available at [www.dic.ac.ir/~EAGLES96/corpusped/corpstpy.html](http://www.dic.ac.ir/~EAGLES96/corpusped/corpstpy.html).

48