

International Corpus of Learner English

Version 2

Sylviane Granger, Estelle Dagneaux, Fanny Meunier,
Magali Paquot (eds.)

UCL PRESSES
UNIVERSITAIRES
 DE LOUVAIN

© Presses universitaires de Louvain, 2009

Registration of copyright: D/2009/9964/3

ISBN : 978-2-87463-143-6

Cover : mikedsign

Printed in Belgium

All rights reserved. No part of this publication may be reproduced, adapted or translated, in any form or by any means, in any country, without the prior permission of Presses universitaires de Louvain.

Distribution : www.i6doc.com, online university publishers

Available on order from bookshops or at

CIACO University Distributors

Grand-Place, 7

1348 Louvain-la-Neuve, Belgium

Tel. 32 10 47 33 78

Fax 32 10 45 73 50

duc@ciaco.com

PREFACE TO FIRST EDITION

Although it has not been around for long, the learner corpus field has already opened up exciting avenues and yielded some interesting results. However, major developments in the field are seriously hindered by the lack of availability of learner corpora. It is our hope that the release of the *International Corpus of Learner English* will mark a new stage in the evolution of EFL research. The corpus, which is the result of over 10 years of international collaboration between a large number of universities, contains 2.5 million words of English written by learners from 11 different mother tongue backgrounds. It is now being released in CD-ROM format with a user-friendly interface which allows researchers to compile their own tailor-made corpora on the basis of a set of predefined criteria. It is accompanied by a handbook which contains a comprehensive description of the corpus, a detailed user manual and brief articles describing the status of English in the countries of origin of the learners. We hope that it will be useful both to SLA theoreticians and EFL practitioners. We are convinced that it can contribute to giving theories of second language acquisition a more solid empirical foundation and lead to the production of more learner-aware pedagogical material, designed for advanced EFL learners in general or focused on the needs of one specific national learner population.

Estelle Dagneaux
Sylviane Granger
Fanny Meunier

December 2002

PREFACE TO SECOND EDITION

The first edition of the *International Corpus of Learner English (ICLE)*, published in 2002, has been very well received internationally. It has been used in a wide range of research projects, served as the basis for many MA and PhD theses and generally played a key role in promoting the emerging field of learner corpus research. The second version differs from the first not only in terms of the higher amount and greater diversity of the learner data included, but also by the improved functionalities it affords. The current version is much more than a simple database of texts: it contains a built-in concordancer which allows for both simple and more complex searches. Another major change is that the rich learner profile information attached to each text is now linked to the search results. In developing the new version, we have taken into account the abundant feedback we have received over the years from scholars using the corpus. We hope the new version will go some way to meeting their requirements. The *ICLE* is a dynamic project: the corpus is regularly updated as new varieties of learner English are added, so there is little doubt that there will be a third release. We therefore encourage users to contact us with their feedback and suggest additional features for future editions.

Estelle Dagneaux
Sylviane Granger
Fanny Meunier
Magali Paquot

February 2009

ACKNOWLEDGEMENTS

First and foremost, we wish to express our deep gratitude to the *ICLE* national coordinators for their relentless efforts in collecting data and the wonderful collaboration we have enjoyed over so many years. A long-term international project such as *ICLE* does not only result in scientific output. It also creates bonds of friendship which will outlast the life of the project. Our thanks also go to the members of the *Centre for English Corpus Linguistics (CECL)*, not only for their precious collaboration at various stages of the project but also for their unfailing enthusiasm. We have also benefited from the dedication and conscientious efforts of numerous students who have contributed to the many painstaking yet necessary tasks that corpus collection entails. Special thanks are due to the beta-testers of *ICLEv2* – Gaëtanelle Gilquin, Przemek Kaszubski and Signe Oksefjell Ebeling – who have provided many useful suggestions for improvement of the interface. We also wish to express our profound gratitude to Paul Rayson, who kindly agreed to POS-tag all the *ICLE* files with the CLAWS7 tagger. Last but not least, we are deeply indebted to the many – well over four thousand – learners who have agreed to contribute their essays to the corpus.

The *ICLE* project would never have got off the ground without the initial and continuing encouragement from three major figureheads in the corpus linguistics field: Jan Aarts, the late Sidney Greenbaum and Geoffrey Leech. By taking a keen interest in the project at a time when the very notion of learner corpus did not even exist, they gave us the necessary confidence to move forward. We gratefully dedicate this volume to them.

We are also greatly indebted to the different funding bodies which have supported the project over the years. The University of Louvain Research Fund played a crucial role in providing us with the initial resources to launch the project. We have also benefited from the generous support of the Belgian National Scientific Research Fund at various stages of the project and from a special research grant from the Walloon Region.

Last but not least, we wish to express warm thanks to Cédric Fairon, Director of the *Centre de Traitement Automatique du Langage*

(*CENTAL*) of the University of Louvain, who welcomed the idea of joining forces on the production of the CD-ROM and helped us to bring the project to a successful conclusion. Special thanks are due to Claude Devis from the *CENTAL* for the wonderful job he did in setting up the *ICLE* database and developing a user-friendly interface for the CD-ROM as well as for his unfailing kindness and patience.

TABLE OF CONTENTS

Preface to first edition	i
Preface to second edition	ii
Acknowledgements	iii
Table of contents	v
List of tables and figures	ix
List of abbreviations	xi
0. Introduction	1
I. Description of the corpus	3
1. <i>ICLE</i> design criteria	3
2. Task variables	4
3. Learner variables	7
3.1. Six clear-cut variables	7
3.2. Two fuzzy variables	10
4. Markup and linguistic annotation	12
4.1. Minimal markup	12
4.2. Built-in concordancer	13
II. Project teams	17
1. Coordinating team	17
2. IT team	17
3. National teams	17
III. ICLEv2 corpus breakdown	25
1. General breakdown	24
2. Breakdown per national subcorpus	25
2.1. Bulgarian	27
2.2. Chinese	27
2.3. Czech	28

2.4. Dutch	28
2.5. Finnish	29
2.6. French	30
2.7. German	31
2.8. Italian	32
2.9. Japanese	33
2.10. Norwegian	34
2.11. Polish	35
2.12. Russian	35
2.13. Spanish	36
2.14. Swedish	37
2.15. Turkish	37
2.16. Tswana	38
IV. ICLE-based research	39
1. Characteristics of <i>ICLE</i> data	39
2. Methodology	40
2.1. Contrastive Interlanguage Analysis	40
2.2. Computer-aided Error Analysis	43
3. Data interpretation	44
V. <i>ICLEv2</i> User Manual	51
1. Introduction	51
2. Licence agreement	52
3. Software installation	56
3.1. Individual licence	56
3.1.1. Installation procedure	56
3.1.2. Launching the software	57
3.1.3. Uninstalling the software	57
3.1.4. Location of the data	57
3.2. Multiple-user licence	58
3.2.1. Installation procedure	58
3.2.2. Launching the software	59
3.2.3. Uninstalling the software	59

3.2.4. Location of the data	59
4. Selecting and querying the corpus	60
4.1. The REQUEST window	60
4.1.1. Navigating the request window	60
4.1.2. Selecting the corpus	60
4.1.3. Keying in a linguistic query	63
a. Single words	64
b. Several consecutive words	64
c. Series of words	64
d. Multiword units (MWUs) or compound lexical entries	65
e. Abbreviations	65
f. Part-of-speech tags and combinations of lexical units and POS-tags	68
g. Morphological filters	73
4.1.4. Note on the composition of the request screens	74
4.1.4.1. Alphanumerical fields	74
4.1.4.2. Numerical fields	75
4.1.4.3. Alphabetical fields	75
4.1.4.4. Variables with multiple options	75
4.1.5. Resetting the selection	76
4.1.6. Functions available on the main command bar	76
4.2. The ZOOM-LIST window	77
4.3. Submitting the query	78
5. The RESULT windows	78
5.1. Two types of result windows	78
5.2. Deselecting profiles	80
5.3. Key terms: Sub-Corpus, Selected Corpus and Result Selected Corpus	81
5.4. Functions available on the result window	83
5.4.1. 'Grid view' and 'form view' of the selected profiles	83
5.4.2. Sorting the data in the Result window (grid view only)	84
5.4.3. Viewing a text and merging texts into a corpus	85
5.5. Statistics	90
5.5.1. Statistics available for corpus collection only	90
5.5.2. Statistics available for joint corpus collection and linguistic query	95
5.5.2.1. More statistics: texts	95
5.5.2.2. More statistics: occurrences	99

6. Frequently asked questions	102
VI. The Status of English	103
1. The Status of English in Belgium (<i>F. Meunier</i>)	104
2. The Status of English in Bulgaria (<i>R. Blagoeva</i>)	112
3. The Status of English in the Czech Republic (<i>V. Minovska</i>)	118
4. The Status of English in Finland (<i>H. Ringbom</i>)	124
5. The Status of English in Germany (<i>G. Lorenz</i>)	128
6. The Status of English in Italy (<i>M.T. Prat Zagrebelsky</i>)	137
7. The Status of English in The Netherlands (<i>P. de Haan</i>)	142
8. The Status of English in Poland (<i>P. Kaszubski</i>)	148
9. The Status of English in Russia (<i>N. Gvishiani</i>)	159
10. The Status of English in Spain (<i>J. Neff, E. Dafouz, J.P. Rica, M. Die & R. Prieto</i>)	166
11. The Status of English in Sweden (<i>B. Altenberg</i>)	172
12. The Status of English in Hong Kong (<i>L. Lin</i>)	178
13. The Status of English in Japan (<i>Y. Ikegami & T. Kaneko</i>)	184
14. The Status of English in Norway (<i>S. Johansson</i>)	190
15. The Status of English in South Africa (<i>B. van Rooy</i>)	196
16. The Status of English in Turkey (<i>A. Kilimci & C. Can</i>)	203
Appendix 1: Institution codes	209
Appendix 2: Suggested essay titles	211
Appendix 3: <i>ICLE</i> -based research: Select list of references	212
Appendix 4: List of compound lexical entries in <i>ICLEv2</i>	214

LIST OF TABLES AND FIGURES

<i>Table 1: Proportion of argumentative essays</i>	5
<i>Table 2: Average essay length</i>	6
<i>Table 3: Top ten essay topics</i>	7
<i>Table 4: Learners' age</i>	8
<i>Table 5: Learners' gender</i>	9
<i>Table 6: CEF results – 20 essays per subcorpus</i>	12
<i>Table 7: Distribution of essays/words per subcorpus</i>	25
<i>Table 8: Bulgarian subcorpus</i>	27
<i>Table 9: Chinese subcorpus</i>	27
<i>Table 10: Czech subcorpus</i>	27
<i>Table 11: Dutch subcorpus</i>	28
<i>Table 12: Finnish subcorpus</i>	29
<i>Table 13: French subcorpus</i>	31
<i>Table 14: German subcorpus</i>	31
<i>Table 15: Italian subcorpus</i>	32
<i>Table 16: Japanese subcorpus</i>	33
<i>Table 17: Norwegian subcorpus</i>	34
<i>Table 18: Polish subcorpus</i>	35
<i>Table 19: Russian subcorpus</i>	36
<i>Table 20: Spanish subcorpus</i>	36
<i>Table 21: Swedish subcorpus</i>	37
<i>Table 22: Turkish subcorpus</i>	37
<i>Table 23: Tswana subcorpus</i>	38
<i>Table 24: List of searchable tags in ICLEv2</i>	69
 <i>Figure 1: ICLE task and learner variables</i>	 4
<i>Figure 2: CLAWS output</i>	14
<i>Figure 3: Unitex input</i>	15
<i>Figure 4: Contrastive Interlanguage Analysis</i>	41
<i>Figure 5: Sample of error-tagged text</i>	43
<i>Figure 6: Error editor screen dump</i>	44
<i>Figure 7: The Query screen</i>	60
<i>Figure 8: Query screen – 'Corpus selection 2' header</i>	61
<i>Figure 9: Merged essays</i>	62
<i>Figure 10: Saving the corpus</i>	63
<i>Figure 11: Query screen – 'Linguistic query' header</i>	63
<i>Figure 12: View word lists option</i>	66

<i>Figure 13: Zoom function icon applied to ‘Other foreign language – first’</i>	77
<i>Figure 14: Typical result window for corpus compilation</i>	79
<i>Figure 15: Typical result window for joint corpus compilation and concordancing</i>	80
<i>Figure 16: Deselecting profiles</i>	80
<i>Figure 17: Deselection of profiles and subsequent deselection of concordance lines</i>	81
<i>Figure 18: Grid view version: Sub-Corpus, Selected Corpus and Result Selected Corpus</i>	82
<i>Figure 19: Form view of the result window</i>	83
<i>Figure 20: Request info</i>	84
<i>Figure 21: Sorting columns in ascending or descending order</i>	85
<i>Figure 22: Generate report option (for variables in the profiles)</i>	86
<i>Figure 23: Main report</i>	86
<i>Figure 24: Generate report options from the concordance lines</i>	87
<i>Figure 25: Report viewer: concordance</i>	88
<i>Figure 26: Report viewer: profiles + concordance</i>	89
<i>Figure 27: Sorting options for the result window</i>	90
<i>Figure 28: Result window for corpus selection – focus on statistics</i>	91
<i>Figure 29: Statistics available for corpus selection only</i>	92
<i>Figure 30: Distribution of the ‘Other learner variables’ in corpus collection</i>	93
<i>Figure 31: Pie chart representation of the statistics (corpus collection only)</i>	94
<i>Figure 32: Result window for joint corpus collection and linguistic query – focus on statistics</i>	95
<i>Figure 33: More statistics: texts – for joint corpus selection and linguistic query</i>	96
<i>Figure 34: Display chart function</i>	98
<i>Figure 35: Chart-like output</i>	99
<i>Figure 36: Statistics related to the occurrences of the linguistic query found in the selected corpus</i>	100
<i>Figure 37: Statistics – task variables</i>	101

LIST OF ABBREVIATIONS

CA	Contrastive Analysis
CIA	Contrastive Interlanguage Analysis
CLAWS	Constituent Likelihood Automatic Word-tagging System
EA	Error Analysis
EFL	English as a Foreign Language
ELT	English Language Teaching
ESL	English as a Second Language
FL	Foreign language
ICLEv1	First version of the ICLE CD-ROM and Handbook (Granger, Dagneaux & Meunier eds. 2002)
ICLEv2	Second version of the ICLE CD-ROM and Handbook
L1	Native language
L2	Foreign/Second language
LOCNESS	Louvain Corpus of Native English Essays
POS	Part-of-speech
SLA	Second Language Acquisition
WST	WordSmith Tools (text retrieval program)

List of key terms provided in the statistics windows of ICLEv2

Corpus	the whole <i>ICLEv2</i> corpus. It consists of 6,085 texts and totals 3,753,030 words.
Sub-corpus	the number of profiles and corresponding texts matching the user's corpus compilation criteria before any potential deselection
Selected Corpus	the number of profiles and corresponding texts left after deselection of profiles (if any)
Result Selected Corpus	the number of texts in the Selected Corpus which contain a linguistic query entered by the user
CorpVar	all texts in the Selected Corpus that match one specific variable

0. INTRODUCTION

This combined CD-ROM and handbook represents the culmination of a project that started in 1990. At the time, corpus linguistics was already a well-established linguistic methodology which was showing its full potential in the field of variation studies. Using a combined quantitative/qualitative approach, corpus linguists were providing much more accurate descriptions of varieties of English than had ever been available before. Although most native varieties of English - regional, diachronic, stylistic - benefited from this new corpus approach, the non-native varieties were completely neglected, which seems strange and unjustifiable when one considers the fact that the number of non-native speakers of English in the world far exceeds that of native speakers. The project launched by Sylviane Granger at the Université catholique de Louvain in October 1990 aimed to bridge that gap. Initially, the project focused exclusively on the writing of advanced French-speaking learners of English, but the idea quickly caught on. Other EFL varieties were added and the project became known as the *International Corpus of Learner English (ICLE)*. A first CD-ROM was released in 2002 (Granger, Dagneaux & Meunier eds.). It contained written data produced by learners from 11 different mother tongue backgrounds: Bulgarian, Czech, Dutch, Finnish, French, German, Italian, Polish, Russian, Spanish and Swedish. From the start we were keen to follow Sinclair's (1991) guidance that "a corpus should be as large as possible, and should keep on growing" and data collection continued uninterruptedly after the release of the first CD-ROM – indeed is still continuing today, with several new national subcorpora still under construction.

With 4.5 million words, the current version of the corpus contains over 1 million words more than the first. The geographical spread of the project has also widened: while the first version only contained data from Europe, the second version includes data from China, Japan and South Africa, in addition to two new European subcorpora (Norway and Turkey). Users of *ICLEv2* will thus benefit from a larger empirical basis for their research. More importantly, they will be able to carry out their research in a totally new technological environment. In *ICLEv1*, users used the CD-ROM to select learner texts that matched a set of predefined attributes but had to use other tools to query the data. In *ICLEv2* the linguistic analysis can be carried out without having to leave the CD-ROM's environment. The built-in

concordancer allows for both simple lexical searches and more sophisticated searches using a POS-tagged version of the data. A major advantage of this integration is that the learner profile information, which was detached from the compiled learner corpora in *ICLEv1*, remains available at all times in the second version, allowing for interesting links between search results and learner/task variables.

The body of this handbook is subdivided into six sections. Part I gives a general overview of the corpus, i.e. the design criteria, the task and learner variables, markup and linguistic annotation. Parts II and III introduce the 16 national projects: the teams that have collected the data (Part II) and the national subcorpora (Part III). Part IV provides some guidelines on how to use the corpus to carry out research into learner language. Part V contains a detailed user manual. Part VI provides information on the status of English in the learners' country of origin.

To order: <http://www.i6doc.com/en/collections/cdicle/>

- Single user licence
- Multiple-user licence (2-10)
- Multiple-user licence (11-25)