

CityU Corpus of Essay Drafts of English Language Learners: A Corpus of Textual Revision in Second Language Writing

John Lee¹, Chak Yan Yeung¹, Amir Zeldes², Marc Reznicek³, Anke Lüdeling⁴, and Jonathan Webster¹

¹ City University of Hong Kong, Hong Kong

² Georgetown University, USA

³ Universidad Complutense de Madrid, Spain

⁴ Humboldt-Universität zu Berlin, Germany

Abstract

Learner corpora consist of texts produced by non-native speakers. In addition to these texts, some learner corpora also contain error annotations, which can reveal common errors made by language learners, and provide training material for automatic error correction. We present a novel type of error-annotated learner corpus containing sequences of revised essay drafts written by non-native speakers of English. Sentences in these drafts are annotated with comments by language tutors, and are aligned to sentences in subsequent drafts. We describe the compilation process of our corpus, present its encoding in TEI XML, and report agreement levels on the error annotations. Further, we demonstrate the potential of the corpus to facilitate research on textual revision in L2 writing, by conducting a case study on verb tenses using ANNIS, a corpus search and visualization platform.

1. Introduction

This article presents a learner corpus that consists of sequences of essay drafts, written by language learners and marked by language tutors. This corpus is designed to facilitate research on how language learners revise their writing, and how feedback influences their revision.

Simply put, learner corpora “have all the characteristics commonly attributed to corpora, the only difference being that the data come from language learners” (Granger 2008)¹. Text corpora often not only contain the raw text, but also supply various kinds of linguistic annotation to facilitate research. For example, grammatical annotation, which can include part-of-speech tags and syntactic structures, is common in many text corpora and also available in some learner corpora (e.g., Reznicek et al. 2013).

A different kind of annotation, particular to learner corpora, is error annotation. It may simply indicate an error category (e.g., Nagata et al. 2011), e.g., marking an inappropriate preposition with the category “wrong preposition”. It may also come in the form of a target hypothesis (e.g., Dahlmeier and Ng 2011; Lüdeling et al. 2008; Nguyen and Miyao, 2013)² — i.e. a corrected or reconstructed version of the learner sentence — in which case the appropriate preposition would also be supplied. Error annotation can be exploited not only in the research of Second Language Acquisition and Foreign Language Teaching (Granger 2004; Nesi et al. 2004), but also in

¹ Examples include the *Cambridge Learner Corpus* (Nicholls 2003), the *International Corpus of Learner English* (ICLE) (Granger et al. 2009), the *National University of Singapore Corpus of Learner English* (Dahlmeier et al. 2013), among many others.

² Target hypotheses are costly to produce and often overlooked, but are nevertheless crucial, since any form of error annotation implies a comparison with what the annotator believes the learner was trying to express. Failing to explicitly document error hypotheses can lead to error annotations that are inconsistent and difficult to rationalize. For extensive discussion, see Lüdeling & Hirschmann (to appear).

automated detection and correction of grammatical errors (Lee and Seneff 2008; Dale and Kilgariff 2011).

To date, learner corpora have concentrated solely on the final form of learner texts, i.e., the end result of the learner's language production process. This process, however, is often an iterative one, with cycles of textual revision, either self initiated or guided by various interventions from the teacher (Graham and Perin, 2007). A corpus containing intermediate versions of learner texts, with feedback, would help us better understand the dynamics of this revision process. It can also potentially provide answers to research questions on many related topics, such as:

- Feedback effectiveness: How often do learners respond to feedback? Which kinds of feedback are most likely to lead to changes, and to what extent do they improve the text (e.g., Truscott 1996; Ferris 1997; Russel and Spada 2006)?
- Revision behavior: How do learners revise their texts, with or without feedback? Which mistakes are most resistant to revision (e.g., Bitchener and Ferris 2012)?
- Language teaching methodology: In view of the above, how can we improve teaching and assessment strategies, and the design of writing assistance tools (e.g., Burstein et al. 2004)?

Some preliminary steps toward data-driven research on these questions have been undertaken. For example, through a web-based EFL writing environment, the XWiLL project offers a searchable database of essays written by students with teachers' comments (Wible et al. 2001); however, the impact of the comments on students' revisions cannot be directly traced. One corpus that aims to address this question is the Malmö University-Chalmers Corpus of Academic Writing as a Process (Eriksson et al. 2012). It is expected to include 450 student texts, ranging from undergraduate to PhD levels, along with peer and teacher feedback. Our corpus is comprised of essays by undergraduate learners, but at a much larger scale, with more than 4,000 essays and over 8 million words. These essays are annotated with error categories and comments from language tutors. Further, for each essay, the corpus includes not only its final version, but also its earlier drafts, with sentence and word alignments (see Figure 2). By combining these annotations and alignments, our corpus provides the largest resource to-date that facilitates research on language learners' revision process, and how it is influenced by feedback.

This article discusses the content, construction of and access³ to the corpus. In the next section, we introduce the compilation, annotation and architecture of the corpus. In Section 3, we report on the conversion process of the corpus material from its original HTML format, as blogs in an e-learning environment, to TEI XML. Although the TEI representation is capable of marking up all the information in the corpus, it still requires considerable programming work to gather non-trivial statistics and create sensible visualizations of the corpus. In Section 4, we discuss how we reduced this technical barrier by importing the corpus to ANNIS, a corpus search and visualization platform (Zeldes et al. 2009). Section 5 presents a case study on verb tense errors using ANNIS. Finally, Section 6 concludes with suggestions for future research directions.

2. Corpus Material

The material in this corpus originated with the Language Companion Course (LCC) project at City University of Hong Kong. The project was implemented over seven consecutive semesters, from 2007 to 2010, involving over 4,200 predominantly Chinese students (Webster et al. 2011).

Essays in the corpus were written by students from across 13 disciplines, representing a wide range of subject areas, including humanities and social sciences, natural sciences and engineering, business, law, and creative media (see Table 1); and a variety of essay genres. Science and engineering courses assigned lab reports⁴; business courses often involved case studies; linguistics students presented data analyses; and social science students wrote argumentative essays. Across all disciplines, article summaries were also assigned.

To support this large body of students, the project recruited more than 300 language tutors, including staff members of the university's English Language Center, and students studying TESOL at one of the university's global partner institutions including undergraduates at Brigham Young University, and post-graduates at the University of Sydney and the University of British Columbia. While full details about the tutors are unavailable, it was the case that whereas those from the University of Sydney and the University of British Columbia came from a variety of language backgrounds, those from Brigham Young University were, for the most part, native speakers of English.

2.1 Drafts

We collected the learner texts and comments from an e-learning environment, the Blackboard system (<http://www.blackboard.com>). Using a word processor built into this environment, the students composed and submitted drafts for written assignments. These drafts were saved as blog entries in HTML format. The tutors then provided feedback on language issues by highlighting problematic words and inserting comments into the drafts. Subsequently, students made revisions to their texts. This cycle continued until the students uploaded a final version to be graded by the professor.

In the rest of this article, each submission is considered a *draft*; a sequence of successive drafts, including the final version, will be referred to as *draft #1*, *draft #2*, etc. One such sequence will be called an *essay*. Our corpus contains 4,337 essays. With an average of 2.7 drafts per essay, there are a total of 11,489 drafts. The average length of a draft is 750 words, yielding a corpus with 8,046,291 words, among the largest annotated learner corpora ever constructed. Detailed statistics can be found in Table 2.

2.2 Error categories

Error annotations were created during the revision process by the tutors. Tutors inserted comments into the drafts in one of two ways: First, they were allowed to select an error category from a fixed list, called the *comment bank*. Adopted from the XWiLL project (Wible et al. 2001) but considerably expanded, the comment bank contains a total of 78 categories, each given a numeric code. In Figure 1, for example, the tutor inserted the code “38” in square brackets, which refers to the error “relative pronoun — missing”. We will call this kind of comment an *error category*. Appendix 1 provides the final version of the comment bank and gives example sentences.

We classified each error category into one of three levels: essay level, clause level, or word level. Table 3 shows the most frequently used error categories at each level. At the essay level, most categories deal with issues of coherence with a few categories relating to informal language

⁴ For lab reports, we include only the discussion section since other sections contain many equations, numbers and sentence fragments.

and the essay structure. At the clause level, categories include punctuation errors, incorrect use of conjunctions and incorrect word order. At the word level, most categories deal with grammatical mistakes, including errors concerning word choice and spelling. The most frequent errors, involving articles, noun number, subject-verb agreement and prepositions, are similar to those found in other English learner corpora, particularly those by native speakers of Chinese and Japanese (Lee and Seneff 2008; Han et al. 2006).

2.3 Open-ended comments

As an alternative to the standard error categories, tutors were also allowed to insert custom comments; these will be referred to as *open-ended comments*. According to previous research, students found detailed comments, specific to their work, to be the most important and useful form of feedback (Lipnevich and Smith, 2009). Consistent with this finding, a meta-analysis also concluded that students benefit more from general explanations of a grammatical phenomenon than from identification of specific errors (Biber et al., 2011). In our corpus, the open-ended comments may explain why a highlighted text was problematic, provide revision suggestions, or raise a question for clarification. In Figure 1, instead of simply using the error category “Pronoun – unclear reference”, the tutor chose to insert the comment [Be more explicit ...] to provide a clearer diagnosis of the problem. While most open-ended comments aim at particular words and phrases, they can also address paragraphs or even the entire essay, such as the comment [Nice report! ...] in Figure 1. Such comments are typically placed either at the beginning or end of a draft.

Tutors were more likely to use error categories than open-ended comments; the former accounts for more than 67% of all comments in the corpus. Both kinds of comments become less frequent, however, as students progress in the revision cycle (see Table 4). For draft #1, more than 2 error categories appear every 100 words; in draft #3 and later, the figure drops to 0.16. A similar trend is observed for open-ended comments. This considerable drop corroborates with previous findings that feedback does help students improve the overall quality of their drafts (Paulus 1999).

Discipline	# Essays	Discipline	# Essays
Applied Physics	288	Electronic Engineering	460
Asian and International Studies	118	General Education	172
Biology	618	Law	20
Building Science and Technology Business	249	Linguistics	644
Business	690	Management Sciences	414
Computer Science	148	Social Studies	477
Creative Media	39		


Table 1. The corpus contains a total of 4337 essays from 13 different disciplines.

	Draft #1	Draft #2	Draft #3+
# Sentences per draft	57.24	60.49	40.63
# Sentences per paragraph	4.50	4.38	4.15
# Words per draft	707.58	782.52	682.45
# Words per sentence	12.36	12.94	16.80

Table 2. Average lengths of drafts, paragraphs and sentences in the corpus. Draft #3+ includes all drafts #3 and beyond.

Clause Level

Comment	Relative pronoun - missing
Explanation	You need to link these phrases with relative pronouns
Examples of Wrong Use	The student gave the presentation made some interesting points.
Correct Use	The student <i>who</i> gave the presentation made some interesting points.
External Links	http://owl.english.purdue.edu/owl/resource/645/01/



Comment Bank by City University of Hong Kong is licensed under a Creative Commons Attribution-Noncommercial-No Derivative Works 3.0 Hong Kong License.

4. They include adding wrong amount of solution in order to have a more accurate, the technique needed to be

dramatic increase of phage titres, reaching to the value of 270million p.f.u./ml. This shows the phenomenon of one-step growth for the T4 bacteriophage.

In each sample time, the phage titres represent the total number of phage that is inside the growth tube, including both the infected virus and [01]suspended virus. Because of it [Be more explicit; what do you refer by 'it'], a sample tube [38]has added in CHCL is used as a control. CHCL can killed infected virus, only the suspensions virus will be counted in the overlay method. By this information, the percentage of viral adsorption can be determined.

In this experiment, error is exist [17]as there is several practical mistake[plural] have been noticed. They include adding wrong amount of solution in the dilution process and using the wrong auto pipette in the mixing procedure. In order to have a more accurate, the technique needed [present tense]to be improved.

[Nice report! You have explained your experiment and your results succesfully. However, it is very important for a Lab report to include background information of previous studies in the Introduction stage, as well as to relate these studies findings with your own results in the Discussion stage. In order to validate your results, to need to support them by referencing other results found in some other studies. In addition, it is essential that in the Discussion stage you make reference to what you stated in your Introduction, so as to reject or support your research purpose. Please bear in mind these comments to prepare assignment 1. I'm looking forward to reading your next piece of work. Best,]

Figure 1. An essay draft annotated with error categories (e.g., [38], with explanations in the pop-up window) and open-ended comments (e.g., [Be more explicit ...]).

Essay-level error categories	# Comments
Informal language	1321
Coherence - More Elaboration is Needed	655
Paragraph - new paragraph	516
Coherence - signposting	315
Coherence - missing topic sentence	191
Clause-level error categories	# Comments
Punctuation - missing	2371
Conjunction Missing	1874
Word order	1577
Punctuation - capitalisation	1475
Sentence - New sentence	1345
Word-level error categories	# Comments
Article missing	10280
Delete this (unnecessary)	9109
Noun - countable	7066
Verb - subject-verb agreement	3929
Preposition - wrong use	3718

Table 3. The five most frequently used error categories, at the essay, clause and word-levels.

Comment type	Draft #1	Draft #2	Draft #3+
Open-ended comments (per 100 words)	33534 (1.24)	17597 (0.87)	1304 (0.15)
Error categories (per 100 words)	60875 (2.26)	24341 (1.20)	1379 (0.16)

Table 4. The number of comments in various stages of the revision cycle.

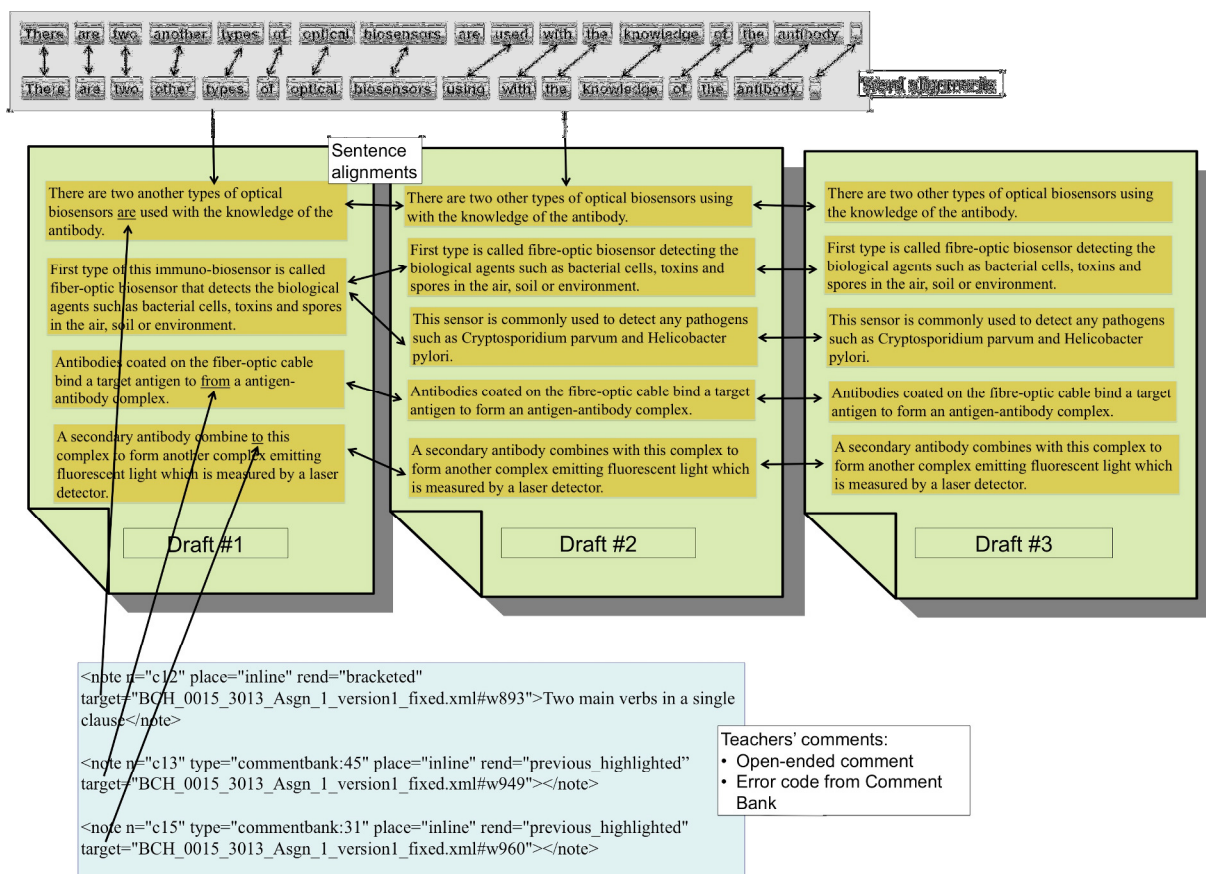


Figure 2. The corpus contains successive drafts of an essay with tutors' comments, as well as sentence and word alignments.

3. Conversion to TEI XML

When the LCC project was conceived, there was no plan to organize the material into a digital corpus. The essay drafts were simply saved as blog entries in HTML format; the comments were marked inline, and not always consistently, in the blogs. Our first task in building the corpus was to convert these blogs into a structured format. Figure 2 graphically depicts the corpus structure.

The Text Encoding Initiative's TEI XML (<http://www.tei-c.org>) is a widely adopted format for text representation and interchange. We chose to encode the corpus according to the current TEI P5 guidelines as this format facilitates further processing by a variety of tools. In our case, we subsequently converted the data for access using ANNIS (see Section 4). The next section explains how each blog entry was stored as two TEI XML files, encoding the essay and the comments respectively. Section 3.2 describes how TEI XML files were generated for sentence and word alignments.

3.1 Drafts

After downloading the drafts as blogs from the Blackboard system, we took the following three steps to convert them to the TEI format.

Automatic Linguistic Annotations

While paragraph breaks can be unambiguously derived from the HTML format, sentence and word boundaries are not explicitly marked. Using Stanford CoreNLP tools (Toutanova and Manning 2000; Toutanova et al. 2003), we split the text into sentences, tokenized the text into separate word forms, and added lemmas and POS tags to the words). After these steps, one TEI XML document was generated for each draft. Paragraphs, sentences and words are enclosed within <p>, <s> and <w> tags respectively, and the parts-of-speech of the words are stored in the ‘type’ attributes of the word tags. Each word and sentence is given a unique id so that it can be referenced from other files.

In order to prevent loss of information, the original appearance of the draft is preserved as much as possible with TEI tags, even if it may not be pertinent to current research questions. For example, we use <hi> to encode formatting styles such as highlighting, bold, underline, strike-through, superscript and subscript. For other special graphical objects that were difficult to capture in text, e.g. pictures and tables, we use appropriate tags like <figure> or <table> without retaining the original graphics.

Mapping comments to text spans

Error categories and open-ended comments are enclosed in brackets and embedded within the drafts. Each embedded comment addresses a particular text span. The text span concerned was indicated either by the font color or the background color of the words in the blog’s HTML format. When an open-ended comment is attached to the beginning or the end of a draft, it is taken to address the entire draft.

The comments are stored in a separate TEI XML file using the <note> tag. Error categories are stored in the ‘type’ attribute while open-ended comments are stored as the text within the tags. The ‘place’ attribute indicates whether the comment is placed in the middle of the draft (usually aimed at a word or sentence), or at its beginning or end (usually aimed at the entire draft). The ‘target’ attribute stores the text span at which a comment is aimed; the text span can be a word element or a range of word elements in the draft.

After this mapping process, we checked whether the error categories are valid, i.e., applicable to the text span at which they are aimed⁵, concentrating on the four most frequently used categories (see Table 3). To be valid, the “article missing” category must be followed by a noun phrase; the “noun countable” category must comment on a noun; the “verb-subject agreement” category must comment on a verb in simple present tense; and the “preposition wrong use” category must comment on a preposition⁶. These categories were found to be valid 96.9%, 99.3%, 97.6%, and 97.4% of the time. We still need to check whether the text spans corresponding to

⁵ Whether the text span contains the specified error is a separate question that will be addressed in Section 3.3.

⁶ We omitted the “Delete this” category, since it can be applied on any kind of word, and so it is always valid by definition.

these valid error categories indeed contain the specified errors; this will be addressed in Section 3.3.

Title and metadata extraction

Most blogs begin with a header, which contains the essay title and metadata such as dates, course codes, grade, assignment and draft numbers, as well as the names and IDs of the student and tutor, which are anonymized. These metadata can facilitate studies on longitudinal improvement, i.e., whether and how a student improved his/her writing through the semester. We extracted the title from the beginning of the header; for ambiguous cases, we compared the extracted title with its counterparts in other drafts of the same essay. Table 5 shows how the title and metadata are stored in components of the <teiHeader>.

Components	Tags	Information
<titleStmt>	<title>	Title of the essay
	<author>	Anonymous student ID
	<editor>	Anonymous language tutor ID
<editionStmt>	<edition>	Draft number (e.g., #1, #2 or #3)
<publicationStmt>	<date>	Date of the assignment
	<idno type="semester">	Year and semester name (e.g., A or B)
	<idno type="course_code">	Course code
	<idno type="assignment_no">	Assignment number, in courses with multiple assignments
	<idno type="grade">	Grade of the assignment (final drafts only)

Table 5. Components of <teiHeader> and the tags and information stored therein.

3.2 Sentence and Word Alignments

To study the revision process, it is imperative to examine how an original sentence in an older draft was edited to form new sentence(s) in the next draft. We automatically obtained sentence and word alignments between drafts and included them in the corpus.

Sentence alignment

This task has been studied in the context of translation of revised documents (Shemtov 1993). Similar to the micro-alignment step in (Barzilay and Elhadad 2003), we used the cosine measure as the lexical similarity metric, and also allow allowed sentence insertion, deletion, merge (two-to-one), and split (one-to-two) alignment. For each consecutive pair of drafts (e.g., drafts #1 and #2, or #2 and #3), the globally best alignment was determined using dynamic programming.

Sentence alignment can be ambiguous. Suppose two sentences, at similar positions in both drafts, share a considerable number of words. The first sentence might have been edited into the second, in which case they should be aligned; alternatively, the first sentence might have been simply deleted and the second inserted, in which case they should not be aligned. Our principle was to prefer higher recall of alignments at the risk of lower precision – i.e. to align sentence pairs with relatively low similarity – since it is much easier for the corpus user to discount an

alignment than to recover an unidentified alignment. This policy was enforced by setting a relatively high cost for insertion and deletion, merge and split.

We chose to use the XCES recommendation for sentence alignments. XCES is the XML application of the Corpus Encoding Standard, a widely accepted set of standards for encoding document structures and linguistic annotations in corpus-based work (Ide et al. 2000). For each consecutive pair of drafts, we encode the sentence alignments in a separate TEI file using <link> tags. Each <link> element has three attributes: ‘prev’ and ‘next’ store the two sentence IDs concerned, ‘type’ stores the alignment type, which may be ‘identical’, ‘edited’, ‘split’, or ‘merged’. A non-aligned sentence may have the alignment type ‘deleted’ or ‘inserted’. There is no ‘next’ attribute in the former case and no ‘prev’ attribute in the latter case.

To evaluate the quality of the automatic sentence alignments, we asked a human judge to manually align the sentences for 14 pairs of drafts. Taking the human alignments as reference, the accuracy of the automatic sentence alignments is 89.8%, measured from the perspective of sentences in the earlier draft.

Word alignment

On the pairs of sentences aligned in the previous step, we further performed word alignment. We obtained word alignments with a tool that calculates the Translation Error Rate, an evaluation metric for machine translation (Snover et al. 2006). This tool generates word alignments as a side product as it calculates the number of word insertions, deletions, substitutions, and shifts between two sentences. Since the “shift” operation allows crossed word alignments, this tool is more suitable for our purposes than most other alternatives.

The words in the sentences are first shifted, or re-ordered, in such a way as to minimize the number of word insertions, deletions, and substitutions. Identical words in the sentences are then aligned. Two non-identical words are considered a substitution (i.e., aligned) if they have similar spelling or have the same lemma. For example, in Figure 2, the aligned pairs “another—other”, and “used—using”, fall into these categories. Otherwise, the words are not aligned, and may be considered an individual insertion or deletion (e.g., “are” in Figure 2).

Similar to sentence alignments, the word alignments are also stored in a separate TEI XML file using the same XCES conventions. Each <link> element has three attributes: ‘prev’ and ‘next’ store the IDs of the two words concerned; and ‘type’ stores the alignment type, which may be ‘identical’, ‘edited’, or ‘shifted’. Similar to sentence alignment, a non-aligned word may have the alignment type ‘deleted’ or ‘inserted’.

3.3 Level of agreement

In order for the corpus to be useful, the reliability of the error annotations is critical. This is usually measured by the level of agreement, i.e., how often two people agree on their error diagnoses – which may involve only error detection or also error correction – on a learner text.

Our evaluation measures how often an expert agrees with error categories annotated in the corpus by the tutors; it is thus an error detection task. Although agreement levels vary depending on error type (Andreu-Andres et al. 2010), perfect agreement is almost never attained for a variety of reasons. Firstly, since learner texts can contain grammatical errors, they often support multiple interpretations (Lüdeling et al. 2008). To judge whether the use of a particular preposition is an error, for example, one must first attempt to reconstruct what the learner “really”

meant. The ambiguous nature of this task is illustrated in studies where subjects were asked to guess which prepositions were originally intended in a set of English sentences. The subjects agreed on the intended prepositions only 76% of the time (Tetreault & Chodorow 2008)⁷; in a similar study, they agreed on the intended article and number for noun phrases only 72% of the time (Lee et al. 2009). To further complicate the matter, the error detection task also demands a judgment on the “acceptability” of a word. Even native speakers often disagree on where to draw the line between a passable word choice and one that ought to be corrected. This difficulty is reflected in a study on the NUCLE corpus, recently used in a shared task for automatic grammar correction (Dahlmeier et al., 2013). When comparing three independent annotations of a sample of 96 essays, the average kappa is 0.39 for grammatical error detection, and 0.55 for error type identification, which correspond to “moderate” and “fair” agreement, respectively (Landis & Koch, 1977). Rosen et al. (2014) reported kappa ranging from 0.16 (“slight” agreement) to 0.88 (“almost perfect” agreement) depending on error type, while Rozovskaya and Roth (2010) attained agreement levels ranging between 56% and 78% on whether a sentence is “correct” or “incorrect”.

Our evaluation focused on the five most frequent error categories (see Table 3). For each category, we randomly selected 200 sentences that contained a text span annotated with that category. We then asked an expert⁸ to decide whether the text span should be revised. Table 6 shows how often the experts agreed with the original annotations by the tutors⁹. The agreement level ranged from 73.9% for “preposition wrong use”, the most challenging category, to 87.1% for “article missing”. These figures corroborate with previous studies in showing error diagnosis on learner text to be highly ambiguous; they also suggest that the reliability of the tutor annotations in our corpus is comparable with existing learner corpora.

Error Category	Agreement level
Article missing	87.1%
Noun countable	78.4%
Delete this	78.9%
Verb – subject-verb agreement	78.8%
Preposition wrong use	73.9%

Table 6. Agreement level between the tutors and the experts in the top five error categories.

4. Access via ANNIS

Although encoded as TEI documents, the corpus still demands considerable programming knowledge and effort¹⁰ on the part of the user to collect meaningful statistics. To reduce the

⁷ This level of disagreement means that evaluation of the precision of error annotations can differ by as much as 10%, depending on the annotator (Tetreault & Chodorow 2008).

⁸ Two experts, both professors of linguistics participated in this evaluation. One was a native speaker of English and the other a near-native speaker who studied in an English-speaking country for 15 years since high school.

⁹ Our evaluation does not estimate the coverage, or recall, of the tutor comments, i.e. the proportion of errors in the learner text that were annotated. Since the tutors were not asked to exhaustively annotate all errors in the text, this figure would not be meaningful.

¹⁰ E.g. using XQuery, a generic query language for XML documents, see <http://www.w3.org/TR/xquery-30/>

technical barrier and to provide a convenient graphical interface to view results, we imported our corpus into the ANNIS system, an open source, browser-based corpus search platform for richly annotated corpora (Zeldes et al. 2009). As an example of its capability, Figure 3 shows a query that returns all sentences with the indefinite article “a” that has been annotated with the error category “article – wrong use” (error category “5”) and revised to “the” in the next draft. We describe the corpus architecture in Section 4.1, then summarize the conversion process from TEI XML to ANNIS in Section 4.2.

4.1 Annotation Layers

Our corpus has various types of annotations that may overlap one another; for example, there can be multiple comments addressing overlapping text spans. Independent annotation layers, encoded in a multilayer architecture, are the most suitable representation, as has been argued before (Reznicek et al. 2013). In such an architecture, any number of types of annotations may be saved in a fashion that prevents one annotation layer from conflicting with another. This allows us to annotate the same category multiple times (e.g. multiple competing part-of-speech annotations), to add different categories to a corpus retroactively without disrupting existing annotations, and to annotate structures that conflict hierarchically, e.g. annotations beginning and ending in the middle of other annotations or discontinuous annotation spans.

The screenshot displays the ANNIS search interface. On the left is the 'Search Form' with the AnnisQL query: `"a" & "the" & CommentBank="5" & #1 ->WordAlign #2 & #1 _#3`. It shows 56 matches in 51 documents. Below is a 'Corpus List' with 'cityu.corpus' selected, containing 11,166 texts and 8,046,291 tokens. The right pane shows search results for the query. It includes a 'Base text' and 'Token Annotations' view. The results are displayed in a table format with columns for words and their POS tags. The first result (14) shows a sentence about a project manager responsible for coordinating. The second result (15) shows a sentence about a design team to resolve engineering. The third result (16) shows a sentence about a team that means. The interface also includes navigation controls and a 'Query Result' tab.

Figure 3. Result of an ANNIS query for sentences with the indefinite article “a” that has been annotated with the error category “article – wrong use” and revised to “the”.

As listed in Table 7, our corpus is represented in ANNIS in nine annotation layers. These layers encode the words in the learner texts, their lemma, POS, formatting style, sentence and paragraph boundaries, and the comments. An example sentence is shown with these layers in Figure 4. Sentence and word alignments are encoded as “relations” between elements in these layers. As listed in Table 8, each relation bears a number of attributes, including the draft number

and the nature of the revision. The interested reader is referred to (Krause & Zeldes, to appear) for further technical detail.

4.2 Conversion to ANNIS

Many different XML formats can be imported into the data model of ANNIS. The most powerful of these in terms of expressivity is PAULA XML (Dipper 2005), which can represent all annotations in our corpus, including annotation spans, parallel alignment on multiple levels and metadata. We converted our TEI documents into PAULA XML, and subsequently imported these into ANNIS using the multi-format, meta-model based converter framework SaltNPepper (Zipser & Romary 2010).

Comment	this should be a phrase, not a clause													
CommentBank	52													
CommentBank	6													
Paragraph	p													
SemesterCommentBank	32													
SemesterCommentBank	1													
Sentence	s													
Style_Render	underline													
Delete	del													
lemma	be	five	of	they	,	each	study	on	particular	aspect	,	include	nature	of service
pos	VBP	CD	IN	PRP	,	DT	NN	IN	JJ	NN	,	VBG	NN	IN NN
tok	are	five	of	them	,	each	study	on	particular	aspect	,	including	nature	of service

Figure 4. A sentence with annotations at the various layers listed in Table 7.

Layer	Description	Layer	Description
tok	Word written by learner	pos	Part-of-speech tag from Stanford tagger
Sentence	Marks start and end of sentence	Comment	Open-ended comment
Paragraph	Marks start and end of paragraph	CommentBank	Error categories
Style_Render	Indicates how the text is rendered (e.g. boldface, underline)	Delete/Insert ¹¹	Whether a word has been deleted or inserted
lemma	Lemma from Stanford tagger		

Table 7. Annotation layers in our corpus.

Attribute	Sentence alignment	Word alignment
From draft	The version number of the draft from which the sentences are aligned.	The version number of the draft the words are aligned from.
To draft	The version number of the draft to which the sentences are aligned.	The version number of the draft the words are aligned to.
Type	‘identical’, ‘replace’, ‘merge’ or ‘split’	‘identical’, ‘replace’ or ‘shift’

Table 8. Annotations of sentence alignment relations and word alignment relations.

¹¹ Although not strictly necessary, this layer improves performance when searching for absence of alignment.

5. Analysis of Textual Revisions: Case Study on Verb Tense

To demonstrate the research potential of the corpus, we present a case study on the influence of feedback on learners' revision behavior regarding verb tense, a common error type in the corpus. Whereas previous studies (e.g., Granger 1999) focus only on the nature of the errors, our corpus enables us to report how often and how these errors are revised, and the impact of feedback on the revision.

This case study focuses on present and past tenses, the most common tenses in the corpus; the four error categories¹² concerned are thus 'verb - present simple' (i.e., revision to present simple is suggested), 'verb – past simple', 'verb – present perfect' and 'verb – past perfect'. There are a total of 2482 comments involving these error categories. We first give an overview of the ANNIS Query Language (AQL) in Section 5.1, showing how it can access, query and visualize relevant materials with a handful of simple queries¹³. We then report an evaluation on the quality of the annotations in Section 5.2, and discuss the results in Section 5.3.

5.1 Queries in ANNIS

Throughout the study, we rely on ANNIS to generate quantitative data. The ANNIS Query Language (AQL) is designed to search for node elements and the edges between them. Roughly speaking, one first specifies the relevant nodes, then states the constraints that must hold between them, and possibly adds metadata restrictions. A node element can be a word, e.g. the query:

```
POS= / (VB|VB[PZ]) /
```

uses a regular expression to find all words tagged as VB, VBP or VBZ, the tags for present-tense verbs for the Stanford tagger, which follows the Penn Treebank tagset (Marcus et al. 1993). Attribute-value pairs can also be used, e.g.

```
CommentBank="85"
```

finds the error category 85, 'verb – past simple', i.e. past simple tense is needed. When specifying multiple elements, the relationship between them must be stated, e.g. both annotations applying to the same position, as in:

```
POS= / (VB|VB[PZ]) / & CommentBank="85" & #1 ==_ #2
```

This query searches for present-tense verbs and the error category 85, and further specifies that the former (#1) covers the same text as the latter (#2), using the operator `==_`. To add the constraint that the present-tense verb was revised to past simple, we add a third search element to find all words tagged as VBD, the tag for past simple verbs. We then use the arrow operator (`->`) and the edge type `WordAlign` to require this past-tense verb (#3) be aligned to #1:

¹² In this study, we do not consider open-ended comments on verb tense errors, since they vary in terms of the explicitness of the feedback, making it difficult to compare their impact. Furthermore, among comments leading to verb tense revision, open-ended comments (16%) are much less frequent than error categories (84%).

¹³ The interested reader is referred to <http://www.sfb632.uni-potsdam.de/annis/> and to (Krause & Zeldes, to appear) for more detail on how the interface can perform sophisticated queries to answer research questions flexibly and without programming skills.

```
POS=/(VB|VB[PZ])/ & CommentBank = "85" & POS="VBD" & #1 __=__ #2 & #1 -
>WordAlign #3
```

In summary, this query searches for cases of a verb in present simple tense which is revised to the past simple in response to the error category ‘verb – past simple’. Some search results are shown in Figure 5.

5.2 Parser evaluation and agreement level

Since our analysis relies on the output of the automatic Stanford POS tagger (Toutanova and Manning 2000; Toutanova et al. 2003), we would like to measure its accuracy on recognizing verb tenses in learner text. We randomly selected 100 sentence pairs involving revised verb tenses, and examined the POS tags assigned by the tagger to the verbs. The accuracy of these tags was 97%. Although automatic syntactic analysis for noisy text is a challenging task (Foster et al., 2008), the tagger seemed capable of correctly analyzing verb tenses in most cases. Our analysis also relies on error annotations by the tutors. Similar to the evaluation in Section 3.3, we used these 100 sentence pairs to measure the level of agreement on verb tense error diagnosis. A human judge annotated these errors, which were then compared with the original annotations of the tutors. The agreement level was 93%, higher than the other categories reported in Section 3.3. Most disagreements appeared to involve the use of the perfect aspect; for example, whether a past/present perfect tense was more appropriate than the past simple, or vice versa.

5.3 Results

As shown in Table 9, it is much more common for students to use the present simple tense where the past was needed, as compared with the reverse direction. This tendency may be influenced by the students’ L1: since Chinese verbs are not inflected, students preferred by analogy the uninflected form in English, which happens to be the present simple.

As for the students’ revision behavior, our corpus shows the feedback uptake rate to be mixed. When asked to change from other tenses to the present simple, students responded at a rate of over 76%; in contrast, when asked to change to the past perfect, they responded in less than 43% of the cases. One explanation could be that they were not as familiar with the past perfect tense as with the present simple. In general, the feedback uptake rate in our corpus is lower than those in the literature. For example, in one study on about 1500 teacher comments, only 14% of the comments were left unaddressed by the students in out-of-class revision (Ferris 1997). A larger study on more than 5700 comments yielded similar conclusions, with only 10% of the comments left unaddressed (Ferris 2006; see also the meta-analysis of a variety of studies in Russel and Spada 2006). Our lower uptake rate may be partially attributed to the fact that the feedback came from tutors rather than teachers.

Even when students did respond to the feedback, it is a separate question whether they improved their writing as a result. For verb tense errors, the rate of improvement again varies according to the tense. When students responded to a suggested revision to the present simple tense, more than two-thirds of their revisions were executed successfully (compare % of changes and % of correct changes in Table 9); when responding to a suggested revision to the past perfect, however, less than a half of their revisions were correct. This discrepancy is consistent with our hypothesis above that the students were unfamiliar with the past perfect. If true, this would point

to the need for giving more detailed feedback for error categories involving complex grammatical constructions.

This case study has investigated only a narrow aspect of the larger research topic of feedback utility, which remains an open question (Truscott 1996; Russel and Spada 2006). When both grammatical feedback and content feedback were given, Fathman & Whalley (1990) reported that all students improved their grammatical accuracy, while 77% also improved the content of their writing. Ferris (1997) found that when ESL students responded to teachers' feedback and made changes, the changes almost always led to overall improvement in their papers. Similar conclusions were made by Ashwell (2000), Ferris and Roberts (2001), Chandler (2003) and Truscott and Hsu (2008). In contrast, Polio and Fleck (1998) found that error correction did not result in any significant improvement in the linguistic accuracy of ESL students. This corpus can potentially contribute further data towards these research questions.

i	to find anything as we think it is just a place	TO VB NN IN PRP VBP PRP VBZ RB DT NN
	to find anything as we thought it was just a place	TO VB NN IN PRP VBD PRP VBD RB DT NN
	⊕ paula	
	⊕ paula text	
i	anything as we think it is just a place with some	NN IN PRP VBP PRP VBZ RB DT NN IN DT
	anything as we thought it was just a place with some	NN IN PRP VBD PRP VBD RB DT NN IN DT
	⊕ paula	
	⊕ paula text	
i	a result , when I try to make my website ,	DT NN , WRB PRP VBP TO VB PRP\$ NN ,
	a result , when I tried to make my website ,	DT NN , WRB PRP VBD TO VB PRP\$ NN ,
	⊕ paula	
	⊕ paula text	

Figure 5. Partial results of an ANNIS query for all present-tense verbs that are changed to past tense in response to a comment to do so.¹⁴

→ Suggested tense ↓ Tense in draft	Present simple	Past simple	Present perfect	Past perfect
Present simple	-	954	89	46
Past simple	345	-	92	21
Past/present perfect	133	131	-	-
Continuous	198	60	12	8
Base form	96	222	49	26
Total	772	1367	242	101
% of changes	76.94%	74.91%	71.07%	42.57%
% of correct changes	56.09%	60.35%	42.56%	20.79%

¹⁴ The query issued was POS=/(VB|VB[PZ])/ & POS="VBD" & CommentBank="85" & #1 ->WordAlign #2 & #1 __ #3. Note that the first sentence is repeated because two verbs were revised.

Table 9. Statistics on the number of times tense-related comments elicited changes and correct changes from the students in their next draft.

6. Conclusion

We presented a corpus containing drafts of a large number of ESL students' essays, together with comments made by language tutors. The corpus incorporates lemma and part-of-speech annotations, and aligns sentences and words from successive drafts, thus showing how students responded to the comments and revised their writing. Currently the largest dataset of its kind, we evaluated the quality of its error annotations, and motivated how it can support research on textual revision in L2 writing.

To provide access to the data, we encoded our corpus in TEI XML format, and further ported it to the ANNIS corpus visualization and search platform. Through a case study, we showed that revision statistics can be retrieved using straightforward queries in the ANNIS Query Language (AQL). This study investigated how ESL students revise verb tense errors and measured their feedback uptake rate. Our analysis indicates that the uptake rate varies according to the tense in question, suggesting that some tenses are more difficult to revise and might warrant more detailed feedback. This was the first study on this topic that is conducted semi-automatically with a query tool on a large-scale corpus.

This corpus can facilitate many directions of research¹⁵. We plan to characterize learners' behavior in textual revision, for example how often sentences are split or merged and at which draft. We would also like to investigate whether and how students revise differently when given or not given feedback; and in the latter case, whether error categories or open-ended comments, and direct or indirect feedback, are more likely to elicit better responses from students.

References

- Andreu Andrés, M. A., Guardiola, A. A., Matarredona, M. B., MacDonald, P., Fleta, B. M. & Pérez Sabater, C. (2010). Analysing EFL learner output in the MiLC Project: An error □it's, but which tag? In Campoy-Cubillo, M. C., Bellés-Fortuño, B. & Gea-Valor, M. L., *Corpus-based approaches to English language teaching*. London: Continuum, 167–179.
- Ashwell, T. (2000). Patterns of Teacher Response to Student Writing in a Multiple-draft Composition Classroom: Is Content Feedback Followed by Form Feedback the Best Method? *Journal of Second Language Writing* 9(3), 227–257.
- Barzilay, R., & Elhadad, N. (2003). Sentence Alignment for Monolingual Comparable Corpora. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*. Sapporo, Japan, 25–32.
- Belz, J. A. and Vyatkina, N. (2005) Learner corpus analysis and the development of L2 pragmatic competence in networked inter-cultural language study: The case of German modal particles. *Canadian Modern Language Review* 62(1):17–48.
- Biber, D., Nekrasova, T., & Horn, B. (2011). The Effectiveness of Feedback for L1-English and L2-Writing Development: A Meta-Analysis. TOEFL iBT Research Report.
- Bitchener, J. & Ferris, D. R. (2012). *Written corrective feedback in Second Language Acquisition and Writing*. New York, NY: Routledge.

¹⁵ Researchers interested in using this corpus are requested to contact the first author to obtain access. The corpus will be made freely available for research purposes.

- Burstein, J., Chodorow, M. & Leacock, C. (2004). Automated Essay Evaluation: The Criterion Online Writing Service. *AI Magazine* 25(3), 27–36.
- Chandler, J. (2003). The Efficacy of Various Kinds of Error Feedback for Improvement in the Accuracy and Fluency of L2 Student Writing. *Journal of Second Language Writing* 12(3), 267–296.
- Dahlmeier, D., & Ng, H. T. (2011). Grammatical Error Correction with Alternating Structure Optimization. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA: ACL, 915–923.
- Dahlmeier, D., Ng, H. T., & Wu, S. M. (2013). Building a large annotated corpus of learner English: The NUS corpus of learner English. In *Proceedings of the Eighth Workshop on Innovative Use of NLP for Building Educational Applications*, 22–31.
- Dale, R., & Kilgariff, A. (2011). Helping Our Own: The HOO 2011 Pilot Shared Task. In *Proceedings of the 13th European Workshop on Natural Language Generation (ENLG)*. Nancy, France, 242–249.
- Dipper, S. (2005). XML-based Stand-off Representation and Exploitation of Multi-Level Linguistic Annotation. In *Proceedings of Berliner XML Tage 2005 (BXML 2005)*. Berlin, Germany, 39–50.
- Eriksson, A., Finnegan, D., Kauppinen, A., Wiktorsson, M., Wärnsby, A., & Withers, P. (2012). MUCH: The Malmö University-Chalmers Corpus of Academic Writing as a Process. In *Proceedings of the 10th Teaching and Language Corpora Conference*.
- Fathman, A. K. & Whalley, E. (1990). Teacher response to student writing: focus on form versus content. In Kroll, B. (ed.), *Second Language Writing: Research Insights for the Classroom*, pp. 178–190.
- Ferris, D. R. (1997). The Influence of Teacher Commentary on Student Revision. *TESOL Quarterly* 31(2), 315–339.
- Ferris, D. R. (2006). Does Error Feedback Help Student Writers? New Evidence on the Short- and Long-term Effects of Written Error Correction. In Hyland, K., & Hyland, F. (eds.), *Feedback in Second Language Writing: Contexts and Issues*. Cambridge: Cambridge University Press, 81–104.
- Ferris, D. R., & Roberts, B. (2001). Error Feedback in L2 Writing Classes: How Explicit Does it Need to be? *Journal of Second Language Writing* 10, 161–184.
- Foster, J., Wagner, J., & van Genabith, J. (2008). Adapting a WSJ-trained parser to grammatically noisy text. In *Proc. ACL*.
- Graham, S. & Perin, D. (2007). A Meta-Analysis of Writing Instruction for Adolescent Students. *Journal of Educational Psychology* 99(3), pp.445–476.
- Granger, S. (1999). *Use of Tenses by Advanced EFL Learners: Evidence from Error-tagged Computer Corpus*. In: Hasselgård, H., *Out of Corpora – Studies in Honour of Stig Johansson*, Rodopi: Amsterdam, Atlanta, 191–202.
- Granger, S. (2004). Computer Learner Corpus Research: Current Status and Future Prospects. *Language and Computers* 23, 123–145.
- Granger, S. (2008). Learner corpora. In Lüdeling, A. and Kyto, M. *Corpus Linguistics: An International Handbook Vol. 1*. Berlin: Mouton de Gruyter.
- Granger, S., Dagneaux, E., Meunier, F., & Paquot, M. (2009). *The International Corpus of Learner English. Version 2. Handbook and CD-ROM*. Louvain-la-Neuve: Presses universitaires de Louvain.

- Han, N.-R., Chodorow, M., & Leacock, C. (2006). Detecting Errors in English Article Usage by Non-Native Speakers. *Natural Language Engineering* 12(2), 115–129.
- Ide, N., Bonhomme, P., & Romary, L. (2000). XCES: An XML-based Encoding Standard for Linguistic Corpora. In *Proceedings of the Second International Language Resources and Evaluation Conference*. Paris: ELRA, 825–830.
- Krause, T. & Zeldes, A. (to appear). ANNIS3: A New Architecture for Generic Corpus Query and Visualization. To appear in *Literary and Linguistic Computing*. <http://llc.oxfordjournals.org/content/early/2014/10/24/llc.fqu057.full.pdf+html>
- Landis, J. R. & Koch, G. G. (1977) The measurement of observer agreement for categorical data. *Biometrics* 33:159–74.
- Lee, J., & Seneff, S. (2008). An Analysis of Grammatical Errors in Nonnative Speech in English. In *Proceedings of the IEEE Workshop on Spoken Language Technology 2008*. , 89–92.
- Lee, J., Tetreault, J., & Chodorow, M. (2009). Human evaluation of article and noun number usage: influences of context and construction variability. In *Proceedings of the Third Linguistic Annotation Workshop*, p.60–63.
- Lee, J., Yeung, C. Y., Zeldes, A., & Lüdeling, A. (In preparation). Textual Revision and Teacher Feedback in Second Language Writing.
- Lipnevich, A. A. & Smith, J. K. (2009). "I really need feedback to learn:" students' perspectives on the effectiveness of the differential feedback messages. *Educational Assessment, Evaluation and Accountability* 21(4), pp 347-367
- Lüdeling, A., Doolittle, S., Hirschmann, H., Schmidt, K., & Walter, M. (2008). Das Lernerkorpus Falko. *Deutsch als Fremdsprache* 2, 67–73.
- Lüdeling, A., & Hirschmann, H. (to appear). Error Annotation. In Granger, S., Gilquin, G., & Meunier, F. (eds.), *The Cambridge Handbook on Learner Corpus Research*. Cambridge: Cambridge University Press.
- Lüdeling, A., Walter, M., Kroymann, E., & Adolphs, P. (2005). Multi-level Error Annotation in Learner Corpora. In *Proceedings of Corpus Linguistics 2005*. Birmingham, UK.
- MacDonald, P., García-Carbonell, A. & Carot-Sierra, J. M. (2013). Computer Learner Corpora: Analysing Interlanguage Errors in Synchronous and Asynchronous Communication. *Language Learning and Technology* 17(2):36–56.
- Marcus, M. P., Santorini, B., & Marcinkiewicz, M. A. (1993). Building a Large Annotated Corpus of English: The Penn Treebank. *Special Issue on Using Large Corpora, Computational Linguistics* 19(2), 313–330.
- Nagata, R., Whittaker, E., & Sheinman, V. (2011). Creating a Manually Error-tagged and Shallow-parsed Learner Corpus. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA: ACL, 1210–1219.
- Nesi, H., Sharpling, G., & Ganobcsik-Williams, L. (2004). Student Papers Across the Curriculum: Designing and Developing a Corpus of British Student Writing. *Computers and Composition* 21(4), 439–450.
- Nguyen, N. L. T. & Miyao, Y. (2013). Alignment-based Annotation of Proofreading Texts toward Professional Writing Assistance. In *Proceedings of the International Joint Conference on Natural Language Processing*, 753–759.
- Nicholls, D. (2003). The Cambridge Learner Corpus: Error coding and analysis for lexicography and ELT. In *Proceedings of the Corpus Linguistics 2003 Conference*.
- Paulus, T. M. (1999). The Effect of Peer and Teacher Feedback on Student Writing. *Journal of Second Language Writing* 8(3), 265–289.

- Polio, C., & Fleck, C. (1998). "If I only had more time:" ESL Learners' Changes in Linguistic Accuracy on Essay Revisions. *Journal of Second Language Writing* 7(1), 43–68.
- Reznicek, M., Lüdeling, A., & Hirschmann, H. (2013). Competing Target Hypotheses in the Falko Corpus: A Flexible Multi-Layer Corpus Architecture. In Díaz-Negrillo, A., Ballier, N., & Thompson, P. (eds.), *Automatic Treatment and Analysis of Learner Corpus Data*. Amsterdam: John Benjamins, 101–124.
- Rosen, A., Hana, J., Stindlova, B., & Feldman, A. (2014). Evaluating and automating the annotation of a learner corpus. *Language Resources and Evaluation* 48:65–92.
- Rozovskaya, A., & Roth D. (2010). Annotating ESL errors: Challenges and rewards. In: *Proceedings of NAACL'10 Workshop on Innovative Use of NLP for Building Educational Applications*.
- Russell, J., & Spada, N. (2006). The Effectiveness of Corrective Feedback for the Acquisition of L2 Grammar: A Meta-Analysis of the Research. In Norris, J., & Ortega, L. (eds.), *Synthesizing Research on Language Learning and Teaching*. (Language Learning & Language Teaching 13.) Amsterdam and Philadelphia: John Benjamins, 133–164.
- Shemtov, H. (1993). Text Alignment in a Tool for Translating Revised Documents. In *Proceedings of the Sixth conference on European chapter of the Association for Computational Linguistics (EACL-93)*. Stroudsburg, PA: ACL, 449–453.
- Snover, M., Dorr, B., Schwartz, R., Micciulla, L., & Makhoul, J. (2006). A Study of Translation Edit Rate with Targeted Human Annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas*. Cambridge, MA, 223–231.
- Tetreault, J. R., & Chodorow, M. (2008). Native judgments of non-native usage: Experiments in preposition error detection. In *Proceedings of the Workshop on Human Judgements in Computational Linguistics*, 24–32.
- Toutanova, K., Klein, D., Manning, C. D., & Singer, Y. (2003). Feature-Rich Part-of-Speech Tagging with a Cyclic Dependency Network. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology (NAACL-HLT 2003)*. Stroudsburg, PA: ACL, 252–259.
- Toutanova, K., & Manning, C. D. (2000). Enriching the Knowledge Sources Used in a Maximum Entropy Part-of-Speech Tagger. In *Proceedings of the 2000 Joint SIGDAT Conference on Empirical methods in Natural Language Processing and Very Large Corpora*. Hong Kong, 63–70.
- Truscott, J. (1996). The Case against Grammar Correction in L2 Writing Classes. *Language Learning* 46(2), 327–369.
- Truscott, J., & Hsu, A. Y.-p. (2008). Error Correction, Revision, and Learning. *Journal of Second Language Writing* 17(4), 292–305.
- Webster, J., Chan, A., & Lee, J. (2011). Introducing an Online Language Learning Environment and its Corpus of Tertiary Student Writing. *Asia Pacific World* 2(2), 44–65.
- Wible, D., Kuo, C.-H., Chien, F.-L., Liu, A., & Tsao, N.-L. (2001). A Web-Based EFL Writing Environment: Integrating Information for Learners, Teachers, and Researchers. *Computers and Education* 37(3-4), 297–315.
- Zeldes, A., Ritz, J., Lüdeling, A., & Chiarcos, C. (2009). ANNIS: A Search Tool for Multi-Layer Annotated Corpora. In *Proceedings of Corpus Linguistics 2009*. Liverpool, UK.
- Zipser, F., & Romary, L. (2010). A Model Oriented Approach to the Mapping of Annotation Formats using Standards. In *Proceedings of the Workshop on Language Resource and Language Technology Standards, LREC-2010*. Valletta, Malta, 7–18.

Appendix 1: Error Categories¹⁶

Complete list of the error categories used in our corpus, with example sentences. The text span addressed by the error category is enclosed in square brackets. For some of the categories, we provide an explanation rather than an example because of space constraints.

Comment	Example/Explanation
<i>Word level</i>	
Adjective comparative / superlative form	The longer I look at the sky, the more free I feel.
Article – unnecessary	[The] oxygen is essential for health.
Article - wrong use	I bought a dress yesterday, but it was the wrong size, so I took [a] dress back to the shop.
Article missing	Please open [] window for me.
Delete this (unnecessary)	The manager must choose the best way [between the solutions] to solve it.
Informal language - Informal language - personal pronouns	You should not reserve seats in the library.
Missing part of the verb-to BE	They [] against this proposal.
Modal – missing	We [] study hard to pass the exam.
Modal - wrong use	He [would] change his job next year.
Noun – countable	Most of the [computer] are new.
Noun – gerund	[Swim] is my favorite sport
Noun – missing	It is true that students always enjoy [].
Noun - uncountable	There is a lot of informations available on this topic.
Part of Speech - incorrect Use - Part of speech – noun needed - Part of speech – verb needed - Part of speech – adjective needed - Part of speech – adverb needed	[Convenient] is very important. The soldiers [defense] the city. It was [convenience]. She sings [beautiful].
Phrasal verb	Please [fill] the questionnaire.
Preposition - missing	To prevent moisture [] entering the plate
Preposition - unnecessary	The video discussed [about] the importance of temperature on the growth of plants.
Preposition - wrong use	I have to finish this work [until] tomorrow
Pronoun - agreement between demonstrative pronoun and nouns	[This] three carbon molecules undergo different reactions.
Pronoun - unclear reference	The fungus secretes an enzyme. However, [it] denatures at a low temperature.
Pronoun - wrong use	Everybody except [she] was sick.
Pronoun missing	They will not come because [] have other plans for the weekend.
Spelling	Fungi can turn a place into a [dessert].
The Genitive	Everyone enjoys their [holiday of two weeks].
Verb - active voice	This crime [was happened] very often.
Verb - bare infinitive	Our teacher made us [to learn] many new words.
Verb - gerund needed	We spent the whole morning [work] on the project.
Verb - intransitive	I [heard] very hard but I still didn't understand.

¹⁶ Due to revisions over the course of the LCC project, the comment bank differed slightly for each semester; in particular, a few categories were annotated at different levels of granularity. For example, “Verb needed”, “Noun needed”, “Adjective needed”, and “Adverb needed” from one semester are subsumed by “Part-of-speech incorrect” from another semester. The more fine-grained categories are considered subcategories in our corpus.

Verb - missing	When I [] up that morning, it was still dark.
Verb - participles	I am [interest] in football.
Verb - passive voice	The results [calculate].
Verb - past perfect	After I [gave] in my work, I was told the rules of the assignment had changed.
Verb - past simple	I [was needing] a great deal of money.
Verb - present perfect	Even though recent studies [revealed] that E.coli is not a comprehensive indicator, this does not affect scientific interest.
Verb - present simple	I [am thinking] he is a good man.
Verb - simple future	If you turn the key to the left two times, the door [opens].
Verb - subject-verb agreement	The secretary to the teachers [were] very smart.
Verb - to-infinitive	I want [learn] a lot of new vocabulary.
Verb - transitive	Our teacher [said] us a long story.
Word choice - Word choice – collocation - Word choice - level of formality	Then I shall [promote to] a university course. The [pics] of the construction site are enclosed.
<i>Clause level</i>	
Conjunction missing OR wrong use - Conjunction - wrong use - Conjunction missing	She drives very fast [so that] she gets a speeding ticket. Fungi help the absorption of nutrients and water for plants, [] they can increase crop production.
Informal language - Informal language - rhetorical	Why were the results faulty?
Punctuation - capitalisation	The two fungi are [pyncoporous] sp. and [cladosporium] sp..
Punctuation - missing	I am working on two projects [] namely nodal and mesh analyses.
Punctuation - wrong use	Since it is raining[.] We will not go hiking today.
Question - do support	Where [] you live?
Relative pronoun - missing	The student gave the presentation [] made some interesting points.
Relative pronoun - unclear reference	The father of John Smith [who] was very young when he became a senator was also quite rich.
Relative pronoun - wrong use	Hong Kong, [that] has a lot of restaurants, offers many different kinds of food.
Sentence - fragment	Fungi helping plants grow.
Sentence - new sentence	She went to buy a new hat[,] it was difficult to find one.
Subject - dangling modifier	Being a prestigious customer of our bank, we are pleased to offer you a 30-day interest free loan.
Subject - dummy subject	[] Rains every day here.
Subject missing	They don't want to pay higher taxes, but [] forces them.
Word order	Not only [this company can] draft a proposal, but also help with promotional activities.
<i>Essay level</i>	
Coherence - drawing a parallel between clauses	The more we know about yoga, we know a lot about the benefits it can bring to us.
Coherence - introductory paragraph missing	<Explanation: a formal essay should have an Introductory paragraph. It is useful if some background information of the topic, the overall viewpoint of the writer and the scope are given in the Introductory paragraph>
Coherence - logical sequence	<Explanation: The student should organise the sequence of the examples more logically.>
Coherence - mismatch between topic sentences and illustrations	<Explanation: The focus of the paragraph is not supported by the illustrations.>

Coherence - missing background information in the introductory paragraph	<Explanation: There is no background information in the introductory paragraph.>
Coherence - missing conclusion	<Explanation: A conclusion is needed to remind the readers of what has been discussed in the previous paragraph.>
Coherence - missing Scope	<Explanation: The main ideas should be outlined in the introductory paragraph to let the readers know what to expect in the upcoming paragraphs.>
Coherence - missing the central focus	<Explanation: There should be one central thought in each paragraph to justify its existence in the essay.>
Coherence - missing thesis statement in the introductory paragraph	<Explanation: The student should insert a thesis statement to state his/her overall opinion or to identify his/her position so that the readers know what to expect in the upcoming paragraphs.>
Coherence - missing topic sentence	<Explanation: A topic sentence is needed here to outline the main idea of this paragraph.>
Coherence - more elaboration is needed	<Explanation: The main idea of each paragraph needs to be supported, explained or illustrated by relevant data, examples, descriptions or explanations to help readers to understand.>
Coherence - signposting	<Explanation: The student can make the relationship between ideas of paragraphs clearer by adding signposting words.>
Coherence - too many focuses in one paragraph	<Explanation: There are so many focuses in this paragraph that the central focus does not stand out.>
Coherence - unclear background information in the introductory paragraph	<Explanation: The background information is not clearly written.>
Coherence - unclear conclusion	<Explanation: A conclusion usually consists of a concluding phrase, a summary of main points, and recommendations. One or more of these elements is missing.>
Coherence - unclear introduction	<Explanation: The student should outline the main ideas in the introductory paragraph to let the readers know what he/she is going to say in the upcoming paragraphs.>
Coherence - unclear scope	<Explanation: The student should give an outline of his/her arguments to prepare the readers.>
Coherence - unclear thesis statement	<Explanation: The thesis statement is not clearly written.>
Coherence - unclear topic sentence	<Explanation: The topic sentence does not outline the main idea of the paragraph.>
Heading – inappropriate	Why fungi can be our friend and foe?
Heading – missing	<Explanation: The student should use keywords to give the focus of the following paragraph(s) / section.>
Illustration	<Explanation: It may be helpful to use an example to illustrate the idea.>
Informal language - Informal language - bullet points - Informal language - contractions - Informal language - headings	There are many reasons contributing to the success of marriage: 1) Communication 2) Patience 3) A caring attitude The government can't resolve this problem quickly. [Analysis of the sentence] We discovered the verbs were foregrounded and this gave a poor balance to the sentence.
Paragraph - new paragraph	<Explanation: One paragraph contains one topic/idea. A new topic/idea needs a new paragraph.>
Reference - Reference - missing or unclear	Britain's most dangerous road is a section of highway linking Lancashire and the Yorkshire Dales.