

Language variation and corpus linguistics

YAMUNA KACHRU*

ABSTRACT: Corpus linguistics deserves serious attention from linguists and applied linguists, since it is of direct relevance to linguistic description, language variation, lexicography, and language education. Linguists tend to be indifferent to corpora, however, as the predominant paradigm in linguistics is based on introspective data, i.e. native speaker intuition. Research has shown that intuitions are not 100 per cent reliable; the notion of 'core' grammar needs to be modified to accommodate the systematic differences across registers at all linguistic levels. Moreover, what linguists perceive as significant principles of linguistic organization may not coincide with their distribution in patterns of use. One goal of applied linguistics is to see what correspondences can be established between the two sets, i.e. the set of underlying principles of linguistic organization and the patterns of use of these principles revealed by analyses of corpora. Regrettably, applied linguists have not embraced corpus linguistics any more enthusiastically than formal linguists. Corpus linguistic analyses have their problems, too. This paper examines a set of data from the lexicon and grammar of world Englishes to suggest that a complete reliance on patterns of use may not solve all the problems of language description, study of variation, language instruction, translation, and lexicography. Furthermore, whereas analyses of corpora are effective in revealing dialect variation, they are not of great use in accounting for diatypic variation, i.e. the permanent characteristics of users of language and recurrent features of their language use, which are crucial for understanding human linguistic behavior.

Move 1
step 1

step2

step3

Move2
step1

INTRODUCTION¹

This paper revisits issues pertaining to corpus linguistics and its application in linguistic description, language variation, lexicography, and language education. Linguistic theory is central to linguistic description, and in its turn, grammar is crucial in the study of language variation and lexicography. The second-order application of linguistics in translation, language education, stylistics, and other areas draw on all of the above. It is evident that linguists tend to be indifferent to corpora, as the predominant paradigm in linguistics, exemplified by transformational generative grammar and its most recent incarnations, is based on introspective data, i.e. native speaker intuition (Chomsky, 1968). Given theoretical linguists' goal of characterizing human linguistic ability, this is understandable. However, the linguistic universals which formal linguists propose have to be evaluated against data from natural languages. Practitioners in corpus linguistics point out that at the level of characterizing the nature of a human language, intuitions are not reliable. The notion of 'core' grammar needs to be amended to accommodate the systematic differences across registers at all linguistic levels (Biber *et al.*, 1994; McEnery and Wilson, 1996). What is true of registers is equally applicable to varieties of a language.

It is true that what linguists perceive as significant principles of linguistic organization may not coincide with their distribution in a particular corpus or a number of corpora, and what grammarians propose as salient patterns of a language may not be the most frequent

*Department of Linguistics, University of Illinois at Urbana-Champaign, 707 S. Mathews Ave., Urbana, IL, 61801, USA. E-mail: ykachru@uiuc.edu

patterns in texts. But these objections do not necessarily mean that the linguist's way of looking at language is less valid; it simply suggests that looking at language as a system and as a medium of communication yield two different sets of results.

APPLIED LINGUISTICS AND CORPUS LINGUISTICS

One aim of applied linguistics is to see what correspondences can be established between the two sets, i.e. the set of underlying principles of linguistic organization and the patterns of use of these principles. Applied linguists have not paid much attention to this task; the debate on the usefulness of corpus linguistics to applied linguistics in general and language education in particular continues (see Seidlhofer, 2003: section 2 for a number of view points on topics related to this debate). A detailed look at the undertaking, however, reveals that the research involved in setting up correlations between grammatical rules and their use in discourse and texts has many applications, especially in descriptions of geographic and social variation, language education, translation, and lexicography, to name just a few areas.

As regards speaker attitudes, beliefs, intentions, etc., linguists interested in pragmatics have described some of the correspondences between linguistic structures and what they signal in terms of the stance of the speaker/writer (e.g. Green, 1989: ch. 6). However, these have not been exploited in applied linguistic projects. For example, Quirk *et al.* (1985: 83–5) categorize both interrogative and negative sentence patterns in terms of non-assertive sentences, in contrast to assertive (declarative) sentences. The justification for this classification is the distribution of sets of items such as the following:

Assertive	<i>some</i>	<i>somebody</i>	<i>something</i>	<i>sometimes</i>
Non-assertive	<i>any</i>	<i>anybody</i>	<i>anything</i>	<i>ever</i>
Negative	<i>no</i>	<i>nobody</i>	<i>nothing</i>	<i>never</i>

Both the non-assertive and the negative forms occur in negative sentences; use of these forms in specific contexts suggests a difference in emphasis. For example, the statement made in utterance (2) is much stronger than that in (1):

- (1) I didn't know anyone there.
- (2) I knew no one there.

Similarly, both *some* and *any* occur in questions, but, as Lakoff (1969) points out, the sentence with *some* reflects a positive speaker attitude, whereas the one with *any* reveals a negative speaker attitude. This may be illustrated with examples such as the following:

- (3) Do you want something to eat before you go to sleep?
- (4) Do you want anything to eat before you go to sleep?

In uttering (3), the speaker hopes for or expects an affirmative answer; in uttering (4), the expectation is of a negative answer. How speakers and writers fully exploit the resources of the grammar and vocabulary in any language is yet to be determined in any detail. Corpus linguistics can contribute to this task if corpora are designed with such analyses in mind. Further examples of such potentialities are discussed below.

Grammatical phenomena

Concerns similar to the one discussed above can be readily exemplified at the level of grammar. English language teaching texts simply take the underlying principles of grammar and present them as though they reflect patterns of use. As a result, one may be faced with a discrepancy between the prominence given to a structural pattern and its actual use as reflected in corpora examined for occurrences of a particular grammatical construction. A good example is the devices for postnominal modification, such as relative clauses and prepositional modifiers: the former are discussed in grammars and teaching texts in detail, whereas the latter are not treated in a similar manner. Should grammars and descriptions of English devote so much more attention and space to the relative clause construction as opposed to the prepositional modifiers? The two structures are exemplified below:

- (5) I have left the books on *the table which is in the hallway*.
- (6) I have left the books on *the table in the hallway*.

A question that should be asked is which is more frequent in a general corpus or in one drawn from a specific genre or genres. According to Biber *et al.* (1994), the prepositional postnominal modifier is the most frequent. A related question is whether the more frequent patterns should be taught first (Biber and Reppen, 2002).

The answers to these questions depend on several factors. However, it is certain that there are two basic conditions that have to be met before teaching texts reflect the distribution of relative clause and prepositional postnominal modification in corpora. First, the relative clause construction is better understood than prepositional postnominal modification, and researchers interested in linguistic description and applied linguistics have to ask themselves why this is the case, and how our understanding of prepositional postnominal modifiers can be improved.

Our understanding of prepositions as a category has to make giant strides in order to meet the latter goal. Although we have substantial evidence to support the claim that English prepositions present a learning problem for most learners, there are no satisfactory texts to teach prepositions in a systematic way. In fact, there is no comprehensive systematic description of prepositions in English grammars. There is still less understanding of the differences between the use of prepositions in the Inner Circle varieties (e.g. American, Australian and British; see Quirk *et al.*, 1985) and between Inner and Outer Circle varieties (Baumgardner, 1996; Bautista, 1997). Some examples of variation in the use of prepositions from Philippine English (PhE) and South Asian English (SAE) are given below:

PhE (Bautista, 1997: 56):

- (7) ... any such venture must be based *from* solid local base.
- (8) ... there are many at this time of the day just *across* this particular library ...

SAE (Bhatia, 1996: 170; Baumgardner, 1987):

- (9) The students ... are trying to escape *out from* this monster of severe disorder.
- (10) Pakistan has no control *to* influence affairs inside Afghanistan.

We need to explore the semantic, pragmatic, and discursal factors that underlie the use of prepositions in Englishes before such uses as those described above can be adequately described. A satisfactory grammatical account of prepositions has to precede a better account of the functions of prepositional postnominal modifiers.

Secondly, we have little data available to show what happens with prepositional post-nominal modification, as opposed to relative clauses, in the learning context. We know, for instance, that learners from certain language backgrounds, e.g. Arabic, Hebrew, and Japanese, have difficulty with the relative clause construction in English (see the discussion in Schachter, 1983), but we have no comparable evidence regarding prepositional postnominal modification. I suspect that once the data becomes available, there will be more discussion of this phenomenon in English textbooks. Effective text material, however, will not be available until the requirement of an adequate description of prepositions is satisfied.

Lexicographical phenomena

Examples from the level of the lexicon add significantly to our appreciation of the task confronting applied linguistics. For instance, the adjectives *hard* and *tough* have very similar morphological properties: they both have comparative forms, *harder* and *tougher*, and accept the verbalizing suffix *-en*, *harden* and *toughen*. However, *hard* has an adverbial form, *hardly*, but *tough* does not have a corresponding *-ly* form. They have similar meanings, i.e. 'resistant to pressure, robust', and both modify human, concrete, and abstract nouns, e.g. *hard bargainer*, *tough boss*; *hard plastic*, *tough steak*; *hard decision*, *tough question*. This grammatical description does not, however, explain why we say *hard feelings*, *hard evidence*, *hard drugs*, *hard cash*, on the one hand, but *tough boss*, *tough skin*, *tough policies*, on the other. The collocation **tough feelings* is as unacceptable as **hard creatures*. For societal use of language, the facts about collocations such as the above are as important as grammatical rules, and corpora can be of immense value in assessing the collocability of items. The responsibility of sorting out collocations motivated by the semantics of lexical items vs. frozen idioms (e.g. *kick the bucket*) still remains, and requires delicate or in-depth linguistic research.

An added complication is introduced by the fact that words mean slightly different things to different speakers, not only across speech communities but within the same speech community. In the case of English, we have to add the speech fellowships of world Englishes within the wider speech community characterized as English speakers. The evidence for this comes from research in several areas, including psychology and computational linguistics, to name just two. Malt *et al.* (1999) investigated the names participants gave to real-world objects, and found that only 2 of the 60 objects presented were given the same name by all 76 native English speaker participants in the study (1999: 242). Reiter and Sripada (2002) report similar findings in technical language use.

In spite of these complications, there is no doubt that corpora are invaluable in lexicographical research, as has been shown by recent studies of lexicons in world Englishes (see e.g. Bautista, 1997; Bautista and Butler, 2000; Butler, 1997; B. Kachru, 1973; 1980; 2006; Pakir, 1992). And yet it is not clear how decisions regarding dictionary listing of senses of an item are to be made – i.e. whether frequency of occurrence or a systematic semantic analysis of the item should be assigned primacy in listing.

This concern with basing listings on corpora analysis can be illustrated with the English lexical item *certain*, also discussed by Biber *et al.* (1994). Assuming that *certain* in the sense of 'sure' is less frequent than in the sense of 'indeterminate', as in *a certain person/book*, etc., in a particular corpus or a variety of corpora, does that difference in frequency mean that the indeterminate sense of *certain* should be listed as the 'core' meaning

of the item? *Certain* in the sense of *sure* enters into some derivational relations (e.g. *uncertain*, *certainty*, *certainly*) and semantic relations (e.g. *certain/uncertain/doubtful*) which it does not in the indeterminate sense. The decision may depend upon the purpose of a particular lexicographical project. If the aim is automated computer processing of texts or teaching English for Special Purposes, the frequency of occurrence has to be taken into account. For linguistic description, however, the systematic semantic analysis of the item will probably take precedence.

Another perspective on dictionary listing can be illustrated with the item *back*. Biber *et al.* (1994) show that the use of *back* in the sense of body part is much rarer than its use as an adverbial, adjective, or verb, which would suggest that the primary meaning of the item is not a body part. However, it is straightforward to argue that *back* designating a body part has something to do with its use as an adverbial, an adjective, and a verb. It would be a pity not to let language learners discover how arbitrary signs such as *back* become non-arbitrary in expressions such as *come back*, *the back door*, *back to the wall*, and *to back into the wall*. As Bolinger says (1980: 24), 'The distinction between the arbitrary and the suggestive is ultimately groundless.'

One other point has to be noted in this regard. As Sampson (1989) has cautioned, existing dictionary listings often exhibit biases against several categories of lexical items, including even common scientific-technical terminologies, negatives, and hyphenated items. In dictionaries of Outer Circle Englishes, specialized vocabularies related to administration, fine arts, law, revenue, sociocultural institutions, etc. are a crucial part of the language of newspapers and other publications. For instance, all early lexicographical compilations in English as used in India were of specialized vocabulary that the British needed to administer the region (see B. Kachru, 2006). Regardless of frequency in a general corpus, items of local importance, technical or not, have to be listed in dictionaries for the dictionaries to be useful to a wide variety of users.

SOCIOLINGUISTIC VARIATION

It has been suggested that a satisfactory descriptive model to account for language variation must distinguish between dialectal variety differentiation and diatypic variation (Oostdijk, 1988). That is, the model must account for permanent characteristics of the users of the language, such as their location in space and time, gender, age, social status, etc., and the recurrent characteristics of the use of language by users, such as field, tenor, and mode of discourse. This is more easily said than done, as is clear to many practitioners in the field of sociolinguistics. It has been impossible to locate any research on such basic features as floor, turn-taking, interruptions, or agreement/disagreement in face-to-face interaction in South or Southeast Asian languages, e.g. in Hindi or Filipino, or the varieties of English used by speakers of these languages (see, however, Valentine, 1995). We need good corpora of spoken material and teams of researchers working on analyses before diatypic variation can be better understood.

Corpus-based research has produced some expected and some unexpected but revealing results in the areas of variation in the use of grammatical and lexical devices. For instance, Collins (1991) focuses on a select subset of modals, those of obligation and necessity, e.g. *must*, *should*, *ought*, *need* and *have (got) to*, and explores their behavior in Australian English (AusE) as compared to American (AmE) and British English (BrE). Each of these modals has two primary meanings: epistemic, signaling the speaker's certainty and

suppositions, and root, indicating obligation, compulsion, or requirement. In their epistemic meanings, *must* and *have (got) to* express a greater degree of conviction than *should* and *ought*. The item *have (got) to* is the main exponent of root obligation in informal speech in AusE (Collins, 1991: 153). Whereas Quirk *et al.* (1985: 225) claim that *must* does not have a negative form (the form *can't* takes the place of *mustn't*), the Australian corpus contains examples of epistemic *mustn't* in conversational data identical to the number of occurrences of epistemic *can't*, with which it is semantically parallel. Since *mustn't* does not occur in written data in Australian corpus, Collins observes: 'Given that linguistic change typically originates in casual spoken genres before spreading to more formal and conservative genres, this distributional pattern suggests that *mustn't* is fairly recent in origin' (1991: 156). It is not surprising that change in context of language use leads to change in language; it is a well-known historical linguistic process, which can be captured by analyses of spoken language corpora.

Current research findings suggest that the corpus-based approach will be effective in coming to grips with dialect and variety differentiation and will deepen our understanding of register and genre variation (see for discussions Aijmer and Stenström, 2004; Biber, 2006; Biber *et al.*, 1998; Biber and Burges, 2000; Conrad and Biber, 2000; and McEnery and Wilson, 1996, among others). However, corpus-based approaches may not meet with the same degree of success in dealing with diatypic variation (Meyer 2002: 17–20). Most difficult to deal with may be tenor, or interpersonal dynamics of discourse, mainly because the various moves in interaction are negotiated as conversations proceed. Further, it is hard to imagine how a corpus-based approach would deal with discourse-based phenomena as opposed to text-based ones. For example, it is not clear how the phenomena of speaker attitudes, beliefs and intentions discussed above in the context of non-assertive and negative forms of determiners such as *some*, *any* and *no* can be dealt with in any corpus-based approach.

There have been some attempts to bring 'subjective' language within the domain of corpus-based research (Wiebe *et al.*, 2004). 'Subjectivity' in this context refers to aspects of language used to express opinions, evaluations, and speculations. Wiebe *et al.* describe studies in which clues of subjectivity are generated from different data sets and tested against other data sets, including low-frequency words (e.g. *monochromatic*), collocations (e.g. *arduous and raucous*), adjectives (e.g. *fascinating, odious*), and verbs (e.g. *stand in awe*). Density of subjectivity clues in surrounding context strongly affects the likelihood of a word's being used subjectively. Such analyses of data may enable researchers to characterize texts as opinion texts and identify subjective discourse segments in texts which in their entirety may not be opinion texts. Even so, this approach may not be applicable to ongoing conversations, debates, negotiations, etc., which need interpretation by a human participant. However, the methodology may be used to sensitize language learners to the phenomenon of subjective language.

CONCLUSION

In all the areas that have been tackled so far, corpus-based linguistic research is as good as the corpora on which it is based, and grammatical or lexical analyses of corpora are as good as the analytical tools, such as grammatical tags or concordances, which are developed to analyze them (Knowles *et al.*, 1996). Furthermore, only limited attempts have been made to carry out semantic and pragmatic analyses and analyses that take into

account sociolinguistic factors. There is plenty of scope for interaction between theoretical linguistics, grammatical description, applied research, and corpus linguistic research.

NOTE

1. I am grateful to Braj Kachru and Cecil L. Nelson for commenting on an earlier version of this paper.

REFERENCES

- Aijmer, Karin, and Altenberg, Bengt (eds.) (1991) *English Corpus Linguistics: Studies in Honor of Jan Svartvik*. London: Longman.
- Aijmer, Karin, and Stenström, Anna-Brita (eds.) (2004) *Discourse Patterns in Spoken and Written Corpora*. Amsterdam: Benjamins.
- Baumgardner, Robert J. (1996) Innovation in Pakistani English political lexis. In *South Asian English: Structure, Use, and Users*. Edited by Robert Baumgardner. Urbana: University of Illinois Press, pp. 174–88.
- Baumgardner, Robert (1987) Utilizing Pakistani newspapers to teach grammar. *World Englishes*, 6(3), 241–52.
- Bautista, MA. Lourdes S. (1997) The lexicon of Philippine English. In *English Is an Asian Language: The Philippine Context*. Edited by M. A. Lourdes and S. Bautista. Sydney: Macquarie Library, pp. 49–72.
- Bautista, MA. Lourdes S., and Butler, Susan (2000) *Anvil—Macquarie Dictionary of Philippine English for High School*. Pasig City: Anvil.
- Bhatia, Vijay K. (1996) Nativization of job applications in South Asia. In *South Asian English: Structure, Use, and Users*. Edited by Robert J. Baumgardner. Urbana: University of Illinois Press, pp. 158–73.
- Biber, Douglas (2006) *University Language: A Corpus-Based Study of Spoken and Written Registers*. Amsterdam: Benjamins.
- Biber, D., and Burges, J. (2000) Historical change in the language use of women and men: gender differences in dramatic dialogue. *Journal of English Linguistics*, 28, 21–37.
- Biber, Douglas, Conrad, Susan, and Reppen, Randi (1994) Corpus-based approaches to issues in applied linguistics. *Applied Linguistics*, 15(2), 169–89.
- Biber, Douglas, Conrad, Susan, and Reppen, Randi (1998) *Corpus Linguistics: Investigating Language Structure and Use*. Cambridge: Cambridge University Press.
- Biber, Douglas, and Reppen, Randi (2002) What does frequency have to do with grammar teaching? *Studies in Second Language Acquisition*, 24(2), 199–208.
- Bolinger, Dwight (1980) *Language: The Loaded Weapon*. London: Longman.
- Butler, Susan (1997) World English in the Asian context: why a dictionary is important. In *World Englishes 2000*. Edited by Larry E. Smith and Michael L. Forman. Honolulu: University of Hawai'i Press, pp. 90–125.
- Chomsky, Noam (1968) *Language and Mind*. New York: Harcourt Brace.
- Collins, Peter (1991) The modals of obligation and necessity in Australian English. In *English Corpus Linguistics: Studies in Honor of Jan Svartvik*. Edited by Karin Aijmer and Bengt Altenberg. London: Longman, pp. 145–65.
- Conrad, Susan, and Biber, Douglas (2000) Adverbial marking of stance in speech and writing. In *Evaluation in Text: Authorial Stance and the Construction of Discourse*. Edited by Susan Hunston and Geoff Thompson. Oxford: Oxford University Press, pp. 56–73.
- Green, Georgia (1989) *Pragmatics and Natural Language Understanding*. 2nd edn 1996. Hillsdale, NJ: Erlbaum.
- Johansson, Stig, and Stenstrom, Anna-Brita (eds.) (1991) *English Computer Corpora: Selected Papers and Research Guide*. Berlin: Mouton de Gruyter.
- Kachru, Braj B. (1973) Toward a lexicon of Indian English. In *Issues in Linguistics: Papers in Honor of Henry and Renée Kahane*. Edited by Braj B. Kachru, Robert B. Lees, Sol Saporta, Angelina Pietrangeli and Yakov Malkiel. Urbana: University of Illinois Press, pp. 352–76. [A revised version in Braj B. Kachru, *The Indianization of English: The English Language in India*. Delhi: Oxford University Press, 1983, pp. 165–89.]
- Kachru, Braj B. (1980) The new Englishes and old dictionaries: directions in lexicographical research on non-native varieties of English. In *Theory and Method in Lexicography: Western and Non-Western Perspectives*. Edited by Ladislav Zgusta. Columbia, SC: Hornbeam Press, pp. 71–101.
- Kachru, Braj B. (2006) English in India: a lexicographical perspective. In *Lexicology 2: An International Handbook on the Nature and Structure of Words and Vocabularies*. Edited by Alan Cruse, Franz Hundsnurscher, Michael Job and Peter Rolf Lutzzeier. Berlin: de Gruyter, pp. 1274–9.
- Kennedy, Graeme D. (1998) *An Introduction to Corpus Linguistics*. London: Longman.
- Knowles, Gerry, Wichmann, Anne, and Alderson, Peter (eds.) (1996) *Working with Speech: Perspectives on Research into the Lancaster/IBM Spoken English Corpus*. London: Longman.
- Lakoff, Robin (1969) Some reasons why there can't be any *some-any* rule. *Language*, 45, 608–15.
- Malt, Barbara, Sloman, Steven, Gennari, Silvia, and Wong, Yuan (1999) Knowing versus naming: similarity and the linguistic categorization of artifacts. *Journal of Memory and Language*, 40, 230–62.

- McEnery, Tony, and Wilson, Andrew (1996) *Corpus Linguistics*. Edinburgh: Edinburgh University Press.
- Meyer, Charles (2002) *English Corpus Linguistics: An Introduction*. Cambridge: Cambridge University Press.
- Oostdijk, N. (1988) A corpus linguistic approach to linguistic variation. *Literary and Linguistic Computing*, **3**(1), 12–25.
- Pakir, Anne (ed.) (1992) *Words in a Cultural Context* Singapore: UniPress.
- Quirk, Randolph, Greenbaum, Sidney, Leech, Geoffrey, and Svartvik, Jan (1985) *A Comprehensive Grammar of the English Language*. London: Longman.
- Reiter, Ehud, and Sripada, Somayajulu (2002) Human variation and lexical choice. *Computational Linguistics*, **28**(4), 545–53.
- Sampson, Geoffrey (1989) How fully does a machine-usable dictionary cover English text? *Literary and Linguistic Computing*, **4**(1), 29–35.
- Schachter, Jacqueline (1983) A new account of language transfer. In *Language Transfer in Language Learning*. Edited by Susan Gass and Larry Selinker. Rowley, MA: Newbury House, pp. 98–111.
- Seidlhofer, Barbara (ed.) (2003) *Controversies in Applied Linguistics*. Oxford: Oxford University Press.
- Svartvik, Jan (ed.) (1992) *Directions in Corpus Linguistics: Proceedings of Nobel Symposium 82*, Stockholm, 4–8 August 1991. Berlin: Mouton de Gruyter.
- Valentine, Tamara (1995) Agreeing and disagreeing in Indian English discourse: implications for language teaching. In *Language and Culture in Multilingual Societies: Viewpoints and Visions*. Edited by Makhan L. Tickoo. Singapore: SEAMEO Regional Language Center, pp. 227–50.
- Wiebe, Janyce, Wilson, Theresa, Bruce, Rebecca, Bell, Matthew, and Martin, Melanie (2004) Learning subjective language. *Computational Linguistics*, **30**(3), 277–308.

(Received 1 October 2007.)

Copyright of *World Englishes* is the property of Blackwell Publishing Limited and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.