

How to use corpus linguistics in sociolinguistics: a case study of modal verb use, age and change over time

Paul Baker and Frazer Heritage

1 Framing sociolinguistics within corpus research

In this chapter we consider how corpora can be used in order to carry out research from a sociolinguistic perspective. Sociolinguistics is a somewhat broad term, with Labov (1972: 183) indicating that it can appear redundant, as all language is social. Despite this, sometimes language can be considered from less “social” perspectives e.g. in terms of a description of how it is structured or to the extent that it resembles other languages. Researchers like Wardhaugh (2005) and Bloome and Green (2002) have identified sociolinguistics as involving consideration of social context and the relationship between language and society. Sociolinguistics can thus involve analysing aspects of language use as they relate to a person’s identity or the community they belong to. It can involve examination of variation between individuals or groups and change over time. In particular, it can concern interactive aspects of language, resulting in corpus-based sociolinguistic studies that have involved concepts like discourse features (Aijmer 2015), politeness (Culpeper and Gillings 2018) and representation of identity (Johnson and Partington 2017). Several studies also have examined where at least two of these concepts overlap (Taylor 2017).

For the purposes of this chapter, we are taking a somewhat narrower focus, considering sociolinguistics from the perspective of speaker identity. We begin by reflecting on why a corpus approach is worth taking in order to answer sociolinguistic questions, followed by a short review of corpus studies of sociolinguistics. We then move on to a case study which is situated as a continuation of research on diachronic change in modal verb usage (e.g. Leech 2002, 2011; Millar 2009). While these studies have used written corpora, we have employed two spoken corpora in order to examine use of the modal verb *may*. The first is the spoken section of the *British National Corpus* (Aston and Burnard 1998), a general corpus of speech and writing collected between 1991 and 1994 (henceforth referred to as the BNC1994), which we compare against the *Spoken British National Corpus* 2014, a second corpus of speech from 2014 (Love *et al.* 2017), referred

to as the BNC2014. Both corpora are relatively similar in terms of containing transcriptions of private conversations among friends and family, although the data collection methods differed slightly.

2 Great challenges bring great benefits

One consideration to be taken into account when using corpora to address sociolinguistic questions is that it tends to be easier to build or find corpora that consist of written texts, particularly if such texts are available in online contexts that already exist in electronic form (as opposed to, say, converting the contents of a handwritten diary). Written corpora *can* be used to examine sociolinguistic variation, particularly if records are kept regarding aspects of the identities of their authors. For example, Cermakova and Farova (2017) have examined variation between male and female authors in corpora of British and Czech fiction. However, many sociolinguistic studies tend to focus on naturally occurring speech, which has traditionally been more problematic for corpus linguists to examine. It can be complicated, expensive and time consuming to collect the large amounts of data required to carry out a meaningful spoken corpus analysis. While the means of recording a conversation has become easier, thanks to the ubiquity of smartphones that can record sound, the resulting audio files still have to be keyed in using a consistent transcription scheme for representing accents as well as para-linguistic and non-linguistic features. Additionally, issues involving ethics, including permission and anonymisation, are, if anything, more salient than they perhaps were 30 years ago. Corpus builders might also struggle to construct a well-balanced corpus that contains speech from a wide range of speaker identities. The BNC2014, for example, contains speech from 671 people and contains a wide range of identity categorisations, with some categories being better populated than others. If we start to combine categories (e.g. counting the possible permutations of age plus gender plus social class), some of the resulting demographic groups will contain very low numbers of speakers. Such studies can be easier if carried out on data collected online, however, e.g. see Subtirelu's (2017) examination of the intersection of gender and race in student evaluations of lecturers on a "rate my professor" website.

With those points considered, though, the benefits of using a corpus (spoken or not) to investigate sociolinguistic issues are great and thus worth engaging with. A large amount of transcribed speech, encoded for different speaker characteristics, presents a huge advantage for analysts, both in terms of claiming representativeness about a particular social group and in terms of allowing existing hypotheses to be explored. Furthermore, this method allows for the identification of linguistic features that may not have been noticed if the analyst had used manual means of identification or was limited to a much smaller dataset. For example, taking a keywords approach to compare different types of women's personal adverts, it was found that women who were seeking relationships with men tended to use the word *me* more often than those seeking relationships with women, resulting in the former group of women making higher use of statements where they positioned themselves as the recipient of actions of the desired other e.g. 'Looking for someone to show me around town' (Baker 2017).

Broadly speaking, corpora have been used in order to examine demographic variation, particularly utilising existing spoken corpora like the BNC1994 see, for example, Schmid's (2003) comparison of lexical items traditionally associated with male and female speech or Rayson *et al.*'s (1997) keyword comparisons of demographic differences. Some studies

have used corpora as part of a larger multi-method study, such as McEnery's (2005) examination of swearing in the BNC1994 which complemented corpus analysis with historical analysis of texts about swearing in seventeenth- and twentieth-century written texts. While lexical and/or grammatical studies have generally been easier to carry out due to the lower requirements in terms of annotation, there have also been (smaller-scale) studies that have considered phonetic or prosodic variation such as Grabe and Post's (2002) examination of stress in different dialects, MacLagan and Hay's (2007) consideration of pronunciation in a corpus of New Zealand speech and Torgersen *et al.*'s (2006) focus on short monophthongs by younger and older speakers in southeast England.

Spoken corpus projects tend to involve data collection over a short period, and it can be difficult to replicate the conditions of an earlier spoken corpus study in order to examine change over time. Therefore, diachronic studies have again tended to involve written corpora, bringing us to a point which this chapter aims to address.

3 Changing frequencies of modal verbs over time

The case study we describe in the following section is based on comparing two spoken corpora to identify change in modal verb use over time. Modal verbs are verbs of necessity or possibility like *would*, *should* or *may* which are often used in conjunction with the base form of another verb (e.g. *might go*). In a diachronic study which compared four British and American reference corpora (known as the Brown Family) containing written texts taken from two time periods, Leech (2002) found that the collective frequencies of 11 modal verbs were lower in the 1990s corpora compared to the 1960s ones. This pattern was more pronounced for the American corpora, which also had lower frequencies of the modal verbs in both time periods when compared to the equivalent British corpora. In particular, verbs which expressed strong modality like *shall*, *need*, *must* and *ought* showed much lower differences in frequency between the two time periods compared to weaker modals like *might*, *will*, *would*, *can* and *could* (with the latter two modals being slightly more frequent in the 1990s American corpus compared to the 1960s one). On the other hand, a group Leech refers to as semi-modals, like *had better*, *want to* and *have (got) to*, had higher frequencies in both 1990s corpora and had higher use overall in the American corpora compared to the British ones. Leech suggests that the patterns reflect a wider trend towards colloquialisation of written language, shown by adoption of features more common to speech (contractions, progressive verbs, questions, genitives, zero relative use), appearing more often in writing. He cautiously discusses other social explanations like democratisation (a tendency to avoid unequal and face-threatening modes of interaction) and Americanisation (the influence of North American habits of expression and behaviour on other nations).

A related corpus study of modal verbs (Millar 2009) analysed the 100-million-word *TIME Magazine Corpus*, containing text from issues of the American *Time* magazine from 1923 onwards. The study found a 22 per cent increase in modals overall between the 1920s and 2000s, although some modals like *shall* had drastically decreased in frequency over time and some, like *would*, showed fluctuation while others like *may*, *can* and *could* had increased. Leech (2011) responded by arguing that the language of *Time* magazine is not representative of all language use, and by including the analysis of additional members of the Brown family of corpora, he shows that modal verb decline is most notable when comparing corpora from the 1960s, 1990s and 2000s. Similarly, his analysis of the *Corpus of Historical American English* (COHA) shows decline in modal

use across the decades of the twentieth century which accelerates from the 1970s. His study notes that *may* is a bone of contention in that it is where his findings disagree most dramatically from Millar (2009). While *may* appears to have decreased over time in the COHA, during the same period, it has increased over time in the TIME corpus. Both Millar and Leech suggest that the rise in frequency of *may* in Time magazine could be due to change in content, style or editorial policy, with a shift towards speculation in reporting and less focus, particularly in magazines, on the past and present – with modal verbs being used to speculate on possible future events in this genre.

In British English, *may* has been associated with politely powerful ways of speaking. Stubbs (1996) describes how the word was used frequently in a court case by a judge during his summing up, in phrases like “you may think/feel/remember” that were addressed to the jury. In some cases, the judge used *may* to signify a course of action he intended to take e.g. “when I sum up I may very well make some comments upon the evidence”, while in others the judge appeared to be politely instructing the jury e.g. “one of the questions you may want to ask yourselves is this”. In the spoken section of the BNC1994, *may* is used more frequently by AB speakers (the highest-earning social class). However, the creation of a second *British National Corpus* (BNC2014) opens up numerous opportunities for the investigation of diachronic variation, as we will demonstrate in the following section.

4 Case study: spoken use of the modal verb *may*

In this chapter, we examine diachronic changes in how different age groups use the modal verb *may*. In order to do this, we use two comparable corpora of speech (the spoken section of the BNC1994 and the BNC2014). The BNC1994 contains 100 million words, although only 12 million of these are transcribed speech. Within the transcribed speech, 7 million words originate from “context-governed” speech – such as language produced in the workplace. The remaining 5 million words derive from private conversations. At present, the BNC2014 contains 11.4 million words of speech, all of which come from private conversations. For comparability, we have only used the 5-million-word sub-corpus of private conversations from the BNC1994 and the whole of the spoken BNC2014. The two corpora are highly comparable and use the same age categorisation system: 0- to 14-year-olds, 15- to 24-year-olds, 25- to 34-year-olds, 35- to 44-year-olds, 45- to 59-year-olds and 60-year-olds and up. Thus, we can study the diachronic change from several perspectives. First, we can compare how the same age groups use language across a single time period, e.g. by comparing all of the age groups in 1994 or all the age groups in 2014. If both time periods show similar patterns, then this would indicate evidence that people’s use of language changes as they get older. But we can also compare speakers of similar age who were born in different time periods, for example, by comparing one set of people who were 0 to 14 years old in the BNC1994 with a second set of people who were 0 to 14 years old in the BNC2014. This kind of analysis allows us to consider possible cohort effects, whether the time period someone was born in will affect their language use. Additionally, we can examine how age cohorts have changed their use of language between the two points of measurement (1994 and 2014). If we assume that the demographically sampled speakers are representative of their cohorts at the time of recording, it would thus be possible to observe how four different cohorts’ use of the modal verb *may* changes as they age. These cohorts are presented in Table 39.1.

Table 39.1 Age cohorts and how old members of those cohorts would have been in the BNC1994 and BNC2014

<i>Cohort birth year</i>	<i>Age range in BNC1994</i>	<i>Age range in BNC2014</i>
People born in the 1950s	35–44	60+
People born in the 1960s	25–34	45–59
People born in the 1970s	15–24	35–44
People born in the 1980s	0–14	25–34

Each of the age groups is representative of people born within decades, rather than across decades. That is to say that the cohorts will stretch from, for example, the birth years of 1950–59 as opposed to, for example, 1954–63. Throughout the rest of this case study, when we refer to age groups, we refer to set age categories (for example, 0- to 14-year-olds). When we refer to decade of birth cohorts, we refer to the groups noted in Table 39.1 (for example, people born in the 1980s).

Any diachronic differences in how members of a given decade of birth cohort use language could reflect two types of change. The first type of change relates to what is seen as socially appropriate uses for a modal verb, such as how the modal *ought* is now seen as archaic. The second change relates to how a birth cohort view their use of the modal verb as they age (or how they believe they should be talking for their age). In other words, diachronic changes could reflect both what is viewed as an appropriate way to use a modal verb in the time period being studied and how age groups apply this to their own language, bearing in mind their age at a given point in time. By taking into account the variables of age group and decade of birth cohort, we can obtain a better picture of variation in the use of *may* among British English speakers at a diachronic level, something which until the availability of the BNC2014 has not been easy to do.

Both corpora are hosted on CQPweb – a freely accessible website hosted at Lancaster University (Hardie 2012). CQPweb has numerous functions which are indispensable to sociolinguists. For example, CQPweb automatically applies the CLAWS tagger to uploaded corpora. The CLAWS tagger annotates the grammatical part of speech for each word and is accurate 96 to 97 per cent of the time (see Fligelstone *et al.* 1997). *May* is polysemous, meaning that if researchers do not specify the part of speech they are interested in, CQPweb will return all forms of *may* – including not only the modal form but also forenames, surnames and references to the month of May. The BNC1994 is tagged using the C5 tag set, and the BNC2014 is tagged using a modified version called the C6 tag set. There are minor differences between the two tag sets (for example, modal verbs are tagged as VM0 in the C5 tag set and as VM in the C6 tag set), but this does not impact on our results. To restrict the results of *may* to only show modal verb uses, we used the search terms “may_VM0” and “may_VM” with the BNC1994 and BNC2014, respectively.

In this study, we use the data disclosed by speakers to investigate how different social groups use language. All the speech in the corpora has been tagged with speaker attributes which include age, gender and social class. These attributes are included in the metadata assigned to each speaker, meaning that analysts can look at how members who hold identity “x” use feature “y”. We use the Distribution function in CQPweb, which allows a researcher to examine how frequently members of a particular social identity

use a specified word. We have also examined concordance lines of occurrences of the features associated with different age groups, leading to a more qualitative analysis based on functional differences.

Our analysis therefore uses two different methods for interpreting the data: the first is to analyse the changes in frequencies at which the modal verb *may* occur. The second is a close reading of extended concordance lines in order to examine if the function of the modal verb has diachronically changed. Within this latter method we also explore the phraseological patterns to examine how speakers use language to achieve the functions associated with the modal verb (requesting permission, hedging and giving permission). These functions were selected because they emerged from a close reading of the concordance lines which contained the modal verb. For both methods, we explore differences by comparing age groups and cohorts. However, in the analysis derived from the second method, we focus on three groups in particular: those who were aged 0 to 14 in the BNC1994, those who were 0 to 14 in the BNC2014 and those who were 25 to 35 in the BNC2014 (members of this group would have been 0 to 14 in the BNC1994).

Before we examine how different age groups use *may*, it is first worth knowing how frequently it is used in the two corpora. Without imposing the restriction of age categories, within the BNC1994, *may_VM0* occurs 637 times across 97 texts (127.03 occurrences per million words). In the BNC2014, *may_VM* occurs 1,365 times across 592 texts (119.50 instances per million words). Thus, the general usage of *may* appears to have slightly decreased across time at a rate of 7.53 occurrences per million words. Bearing this in mind, Figure 39.1 outlines the diachronic differences in the frequency at which *may* (as a modal) occurs across comparable age groups.

The data from Figure 39.1 paint a rather complex picture: Across all six comparable age groups, three (0- to 14-year-olds, 25- to 34-year-olds and 45- to 59-year-olds) appear to use the modal verb *may* more frequently in 2014 than in 1994. In contrast, the other three (15- to 24-year-olds, 35- to 44-year-olds and 60+-year-olds) appear to use *may* as a modal verb less often in 2014 than in 1994. While the frequency at which *may* occurs as a modal in the BNC2014 as a whole has decreased by 7.53 per million words, when we examine how often different groups use it, the amount of change ranges from 0- to 14-year-olds, who use *may* 63.36 times per million words more often in 2014 compared to 1994, to 35- to 44-year-olds who use *may* 38.99 occurrences per million words less in 2014 compared to 1994. These differences in usage compared to the general picture of diachronic changes to the modal form of *may* give weight to why sociolinguistic analyses of corpora are needed. Although examinations of diachronic change in the frequency at which a word is used can provide interesting findings, they do not necessarily reflect the behaviours of different social groups, which can deviate from the perceived diachronic change in a general sense.

Does Figure 39.1 show evidence of an age effect? Considering just the bars for 1994, there does appear to be a shift in frequency which is weakly related to age, with the younger three age groups (people aged under 35) using *may* less than the older three age groups (people aged 35 and over). However, the bars for 2014 show a different pattern, with no perceivable linear trend. Therefore, age differences in *may* do not appear to be consistent across the two points of measurement, indicating that the age effect cannot fully explain use of *may*.

What about the cohort effect? We can take the data from Figure 39.1 and arrange the bars differently to compare how the birth cohorts use *may* across the two time periods (Figure 39.2).

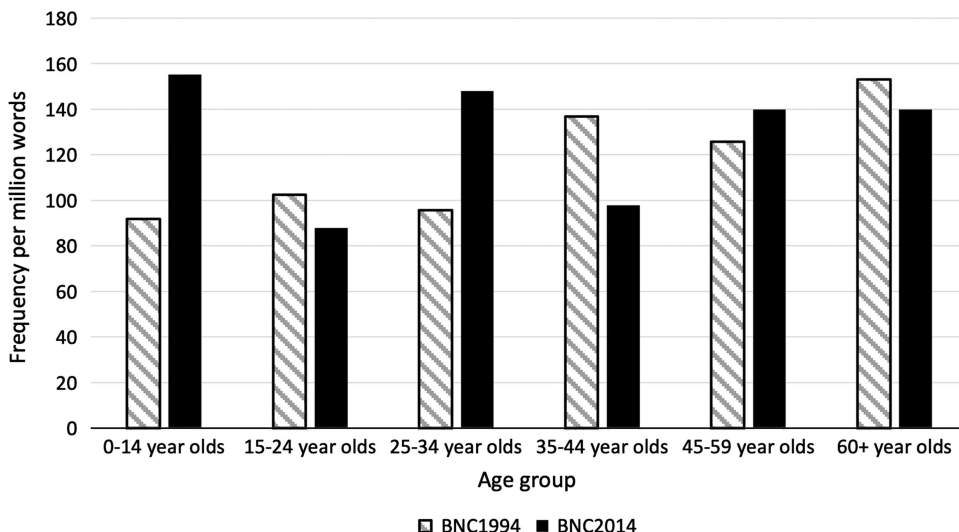


Figure 39.1 Diachronic change in *may* as a modal verb across comparable age groups

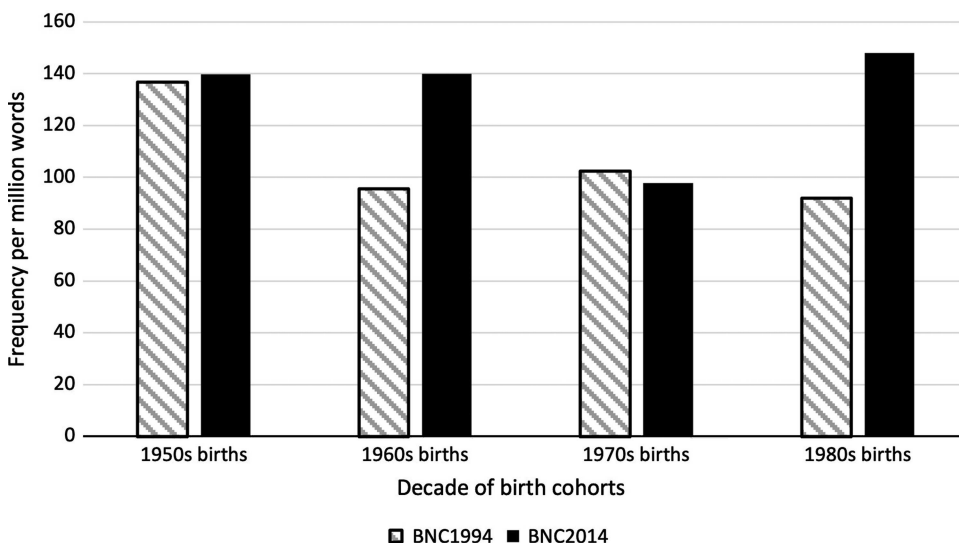


Figure 39.2 Diachronic changes in how frequently decade of birth cohorts used the modal verb *may*

Figure 39.2 indicates a different picture compared to the data presented in Figure 39.1. The 1950s and 1970s births show very little change in their use of *may* between the two time periods when they were recorded. Of course, this should not indicate their use of *may* has been constantly at the same level for all the years between 1994 and 2014, but we can at least say that when 1994 and 2014 are compared, there are only tiny differences in relative frequency for these two age cohorts. While the 1950s births have high use of *may* in both periods, the 1970s births have relatively low use of it.

However, the 1960s and 1980s births follow a different pattern: lower use of *may* in 1994 and higher use in 2014. Given that there are cohort differences which do not match the age group differences, it appears that both the time period a person is living in (whether it is the 1990s or 2010s) and their age at the time play a part in the frequency which they use *may*.

Quantitative findings such as these can provide a useful broad picture of diachronic changes in how a word is used. However, it is useful to combine them with an analysis of context, as this can reveal how a particular word is used in different time periods by different age groups and members of decade of birth cohorts.

In order to gain a clearer understanding of how *may* is used in both corpora, we conducted a close reading of concordance lines. When working with reasonably large frequencies, it can often be useful to examine collocates of a word (see for example Hunston, this volume; Jones, this volume), as a way of down-sampling concordance lines to a manageable amount. However, when we considered occurrences of *may* based on use within particular age groups, we found that its frequency was too low to produce collocates that met traditional significance thresholds (for a discussion of these thresholds see Brezina 2018). Given the limitations of space, we only focus on three groups: those who were 0 to 14 years old in the BNC1994, those who were 0 to 14 years old in the BNC2014 and those who were 25 to 34 years old in the BNC2014. These groups were selected for closer analysis because 0- to 14-year-olds in the BNC1994 and the BNC2014 had the greatest difference in normalised frequencies at which *may* occurs as a modal verb (see Figure 39.1). In other words, 0- to 14-year-olds in the BNC2014 used the modal verb *may* much more than their comparable group from the BNC1994. Furthermore, the group who were 0 to 14 years old in the BNC1994 (who are the same decade of birth cohort as the 25- to 34-year-old group in the BNC2014) had the largest increase in the normalised frequency of use of *may* as modal as they aged (see the 1980s births in Figure 39.2). We thus start by analysing how these age groups used the modal verbs in the BNC1994 compared to the BNC2014 in terms of whether they use it to request permission (e.g. *may I have some milk*), hedge propositions (e.g. *you may want to go*), give permission (e.g. *you may do that*) or other uses (e.g. cases that were unclear due to interruptions). We examined all concordance lines in which *may* was used (as a modal verb) by a speaker within these groups. These quantified frequencies (as a percentage of the raw total for each group) are presented in Figure 39.3.

The figure indicates notable changes – both diachronically across the comparable age group and within the 1980s decade of birth cohort. First, across the comparable age groups (consider the first two sets of columns), the 0- to 14-year-olds in the BNC1994 appear to use *may* as a modal to request permission much more than their 2014 counterparts, who instead use *may* more to hedge a proposition.

We can use Figure 39.3 to focus just on people born in the 1980s by comparing the first and third sets of bars together. In 1994, this age-cohort were children, and as noted earlier, they tended to use *may* to request permission. However, in 2014 they were in their late twenties/early thirties, and use of *may* was used much more often to hedge a proposition.

The middle set of bars in Figure 39.3 thus resembles the final set of bars more than it does the first set of bars, indicating that the children of 2014 are more similar in their use of *may* to older people in the same time period as them, compared to the children of 1994. The figure also indicates some evidence that the function of *may* might be changing over time – from being used to make polite requests to instead hedge propositions.

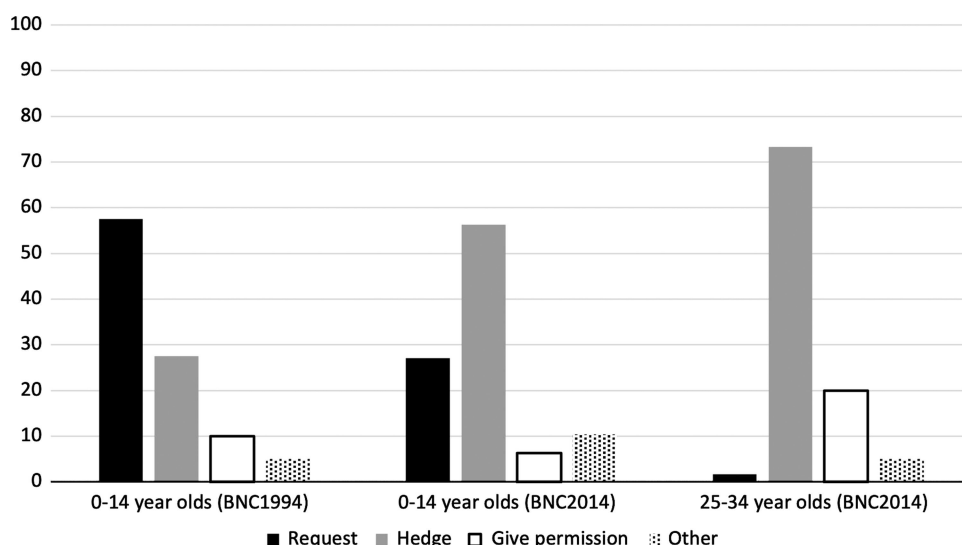


Figure 39.3 Percentages of the functions of *may* in three groups

However, we have only examined three age groups and would need to make a fully study of all the other age groups across both time periods in order to investigate this hypothesis further.

The analysis indicates that there are not only differences in how comparable age groups use language, but there are also differences in how members of the same decade of birth use language. However, while quantifying how frequently these different functions occur may be helpful in painting a picture of diachronic change on a quantitative level, it does not necessarily explain how the modal verb is used at a phraseological level. That is to say, it does not necessarily reveal how *may* is used to achieve these functions. In order to do this, we examined the phraseological patterns around *may* for these groups within the datasets.

Table 39.2 notes the different phraseological lexical units in which *may* was used. In order to ascertain these, we used the “sort” function in CQPweb. This function orders concordance lines by alphabetising co-occurring words in a pre-set slot to the left or right. Although analysts have control of which slot is set for the sort function, the results we present are based on sorting the concordance lines by the co-occurring word one place to the right. We only noted a phraseological pattern if it occurred in at least 10 per cent of the concordance lines for a particular group. This prevented low-frequency patterns being taken as representative usages.

One immediate difference between the language used by 0- to 14-year-olds in the BNC1994 and the BNC2014 is that those in the BNC1994 appear to use *may I* much more regularly than those in the BNC2014. In both datasets, the modal construction of *may I* is used in order to request permission to do something and to request physical items, which could explain the difference in the frequency at which *may* is used to request permission demonstrated in Figure 39.3. However, the 0- to 14-year-olds who were recorded for the BNC2014 use a greater variety of phraseological constructions which incorporate *may*, compared to children of a similar age who were recorded for the BNC1994. Additionally, the phraseological constructions appear to suggest that the

Table 39.2 Common phraseological units containing *may* for different age groups

Age group (Data set)	Raw frequency	Notable phraseological lexical units(raw frequency of occurrence) example (filename)
0- to 14-year-olds (BNC1994)	40	<i>may I</i> (21) <i>may I have my pudding please?</i> (KBW 9002)
0- to 14-year-olds (BNC2014)	48	<i>may be</i> (5) <i>it may be a chocolate goose it may be a giant Easter egg</i> (S8LG 404) <i>may have</i> (6) <i>she touched him and got it and er and she may have been electrocuted I can't remember</i> (SNZS 175) <i>may I</i> (12) <i>may I have some cake?</i> (SCG9 533) <i>may not</i> (8) <i>she thinks she's a witch but she may not be a witch</i> (S46J 244)
25- to 34-year-olds (BNC2014)	240	<i>may as well</i> (24) <i>if there's a little two you can make you may as well try and make it</i> (S968 1083) <i>may be</i> (43) <i>United States somewhere that's bigger where if we get where we may be able to afford property</i> (SVZB 1386) <i>may have</i> (33) <i>he may have been drunk when they were talking about this deal</i> (S4S5 631)

0- to 14-year-olds in the BNC2014 use constructions containing the modal form of *may* to talk about events they remember from the past, request items and hedge ideas (such as seeing if something is a goose or an egg in Table 39.2). This contrasts with the 0- to 14-year-olds from the BNC1994, who mostly only use it to request items. In turn, these phraseological patterns provide some qualitative explanation for the data presented in Figure 39.3.

In contrast to both 0- to 14-year-old groups, the 25- to 34-year-old group in the BNC2014 appears to use *may* in ways which suggest uncertainty about more abstract concepts, such as finances and biological phenomena. Furthermore, the 0- to 14-year-olds in the BNC2014 exclusively use *may have* to refer to past events. In the BNC1994, those who were 0 to 14 only used *may have* three times (so it was not considered a phraseological lexical unit): two of these occurrences were references to past events, and one was a request. Comparatively, the 25- to 34-year-olds in the BNC2014 appear to use *may have* to refer to both past events and potential future events, for example: *we may not need the twelfth month we may have finished all our work and done everything and therefore that would save us like seven hundred and fifty quid in rent* (SBTC 1149). Indeed, for this group there are 7 occurrences (out of 33) which use *may have* to refer to future possibilities, all of these refer to events in the distant future, rather than the immediate future.

When specifically examining the difference between how members of the decade of birth cohort born in the 1980s use language, it can be seen that they no longer retain the frequent use of *may I*. This politeness form, which appears to have been appropriate for 0- to 14-year-olds in 1994, appears to have been discarded later in life. Furthermore, the fact that *may I* is used much less by children born in later generations, indicates that perhaps it is not being taught to children as an appropriate way to make requests.

5 Reflecting and future directions

This case study has demonstrated how corpora can be used to explore not just social variation but also diachronic changes with regard to the use of modal verbs. It is now worth stepping back and considering some of the issues with using corpora in this way. In particular, the study raises issues regarding essentialist approaches to language. Throughout, we have tried to avoid the claim that “x” age group use “y” feature because they are part of the “x” age bracket. To exemplify this problem, one must ask “when someone has a birthday which pushes them into the next age bracket, do they suddenly stop using a linguistic feature?” The likely answer to this is no. In turn, this serves to demonstrate the issues of using statically categorised social groups. Although we have done our best to avoid essentialist claims, the practice of putting people into categories and then counting frequencies of linguistic features does face the risk of appearing essentialist, particularly as we cannot account for every intersection of a speaker’s identity (age, social class, gender), and even when we are able to account for a number, these identities may not even play a role in the language used by that speaker (for a discussion of intersectionality, see Crenshaw 1990). Importantly, to avoid essentialist thinking, a sociolinguistic corpus analysis should first note that differences in frequency are usually not absolute, but more a matter of gradience, with one social group using a feature more than others, but often not having full ownership of that feature. To that end, the qualitative forms of analysis favoured by examining concordance lines allow us to consider differences within a particular social group as well as differences between them. Of course, it is often the case that a corpus search will yield too many concordance lines and in such cases the analyst might decide to analyse a random sample of, say, 100 lines, noting the main trends and deciding to look at further lines if the patterns appear inconclusive. Considering context of usage via an analysis of an expanded concordance line allows us to identify individual speaker variation that can indicate that people within a group can use a word for different purposes (and indeed, the same speaker can vary in the way they use a word). Such considerations will help us to provide a more sophisticated sociolinguistic analysis that goes beyond a table of numbers and claims that group “x” use feature “y”.

Given the limitations of space and the limitations of the corpora (which would need to be even bigger to downscale sample sizes to account for more identities), it has not been fully possible to explore all intersections of identity in our exploration. While it would be interesting to compare these changes to how other social groups use the same feature and how adding other dimensions changes the variation, even larger corpora would be required to do so. However, we should return to a point made earlier in this chapter, regarding the fact that when we divide speakers into groups, the larger the number of groups, the smaller the frequencies. Particularly for speakers aged 0 to 14 in both spoken BNCs, the amount of data to work with is smaller compared to the older speakers. There are 435,286 words of 0- to 14-year-olds’ speech in the BNC1994 and 309,177 words in the BNC2014. Particularly for the latter corpus, word counts for the older age groups are much higher (e.g. 2,777,761 words for people aged 15 to 24, while this was only 795,250 for the same age bracket in the BNC1994). As a result of the smaller amounts of children’s speech, the raw frequencies of *may* are low in both corpora (40 and 48 occurrences), meaning that classification into functional categories results in even lower frequencies, giving less confidence that results are generalisable and would be replicated if carried out on an equivalent sample. Also, the age group 0 to 14 covers an

extremely crucial period of language development, resulting in very marked changes, so care must be taken not to generalise use of *may* as being similar for children aged 2 compared to a teenager of 14. When using spoken corpora, we thus need to consider the extent to which identity categories contain a workable amount of data and are actually useful constructs for our purposes.

Future research that utilises corpus approaches to sociolinguistics may want to focus on using corpus methods on data generated from communities of practice (see Eckert and McConnell-Ginet 2007). By using language from communities that are known to organise around a particular identity, it might be possible to compare and contrast variation which is community specific, and avoid some of the essentialism that naturally occurs when using large reference corpora. Additionally, future research may want to examine both different modal verbs and how they are used within the social groups that are the focus of this chapter. It would be interesting, for example, to examine other modal verbs used to signal politeness and compare this across age groups. Elsewhere, other researchers may want to examine how the modal verb *may* changes across different identities and how and why these identities might play a role in the use of the modal verb *may*. Having noted that children no longer seem to use the *may I* form to make requests, we might want to consider alternative forms that they might use e.g. *give me* or *can I*.

Finally, this study has only been able to answer the question “how do these groups use language?” Future research might elect to attempt to answer the question “why do these groups use these features?” We might want to form hypotheses based on the patterns found e.g. younger children are taught formal, polite forms of language less than used to be the case, which anecdotally, might be worth exploring. Analysing the data itself might also provide clues (although we did not find any cases in either corpus of older people correcting children’s language e.g. “No, if you want something, say *may I have...*”). Instead, a deeper and more nuanced study, which uses interviews, discourse completion tasks, focus groups or media texts to examine attitudes towards these kinds of modal verbs and how people recall their own experience of language use and ideologies as they were growing up, could complement this one, providing further insights into this kind of social variation.

Further reading

- Baker, P. (2010) *Sociolinguistics and Corpus Linguistics*, Edinburgh: Edinburgh University Press. (This book acts as a general primer for a range of ways that corpora can aid sociolinguistics, having chapters on demographic variation, comparing language use across different cultures and examining language change over time, studying transcripts of spoken interactions and identifying attitudes or discourses.)
- Friginal, E. (ed.) (2017) *Studies in Corpus-based Sociolinguistics*, London: Routledge. (This edited collection of 14 chapters from a range of authors is divided into three sections: languages and dialects, social demographics and register characteristics.)
- Friginal, E. and Hardy, J. (2013) *Corpus-based Sociolinguistics: A Guide for Students*, London: Routledge. (This book functions as a practical guide for students who wish to carry out their own studies, containing case studies, discussion questions and activities.)
- Murphy, B. (2010) *Corpus and Sociolinguistics: Investigating Age and Gender in Female Talk*, Amsterdam: John Benjamins. (This monograph involves a detailed analysis of age and gender in a 90,000-word spoken corpus of Irish English, considering features like hedges, vagueness, intensifiers and swearing.)

References

- Aijmer, K. (2015) 'Analysing Discourse Markers in Spoken Corpora: *Actually* as a Case Study' in P. Baker and T. McEnery (eds) *Corpora and Discourse Studies*, London: Palgrave, pp. 88–109.
- Aston, G. and Burnard, L. (1998) *The BNC Handbook: Exploring the British National Corpus with SARA*, Edinburgh: Edinburgh University Press.
- Baker, P. (2017) 'Corpora and Social Demographics', in E. Friginal (ed.) *Studies in Corpus-Based Sociolinguistics*, London: Routledge, pp. 157–77.
- Bloome, D. and Green, J. (2002) 'Directions in the Sociolinguistic Study of Reading', in P. D. Pearson, R. Barr, M. L. Kamil and P. Mosenthal (eds) *Handbook of Reading Research*, Mahwah, NJ: Lawrence Erlbaum, pp. 395–421.
- Brezina, V. (2018) *Statistics in Corpus Linguistics: A Practical Guide*, Cambridge: Cambridge University Press.
- Cermakova, A. and Farova, A. (2017) 'His Eyes Narrowed—Her Eyes Downcast: Contrastive Corpus-Stylistic Analysis of Female and Male Writing', *Linguistica Pragmensia* 28 (2): 7–34.
- Crenshaw, K. (1990) 'Mapping the Margins: Intersectionality, Identity Politics, and Violence against Women of Color', *Stanford Law Review*, 43(6): 1241–99.
- Culpeper, J. and Gillings, M. (2018) 'Politeness Variation in England. A North-South Divide?', in V. Brezina, R., Love and K. Aijmer (eds) *Corpus Approaches to Contemporary British Speech*, London: Routledge, pp. 33–59.
- Eckert, P. and McConnell-Ginet, S. (2007) 'Putting Communities of Practice in their Place', *Gender & Language* 1(1): 27–37.
- Fligelstone, S., Pacey, M. and Rayson, P. (1997) 'How to Generalize the Task of Annotation', in R. Garside, G. Leech, and A. McEnery (eds) *Corpus Annotation: Linguistic Information from Computer Text Corpora*, Longman, London, pp. 122–36.
- Grabe, E. and Post, B. (2002) 'Intonational Variation in English', in B. Bel and I. Marlin (eds) *Proceedings of the Speech Prosody 2002 Conference*, 11–13 April 2002, Aix-en-Provence: Laboratoire Parole et Langue: 343–46.
- Hardie, A. (2012) 'CQPweb - Combining Power, Flexibility and Usability in a Corpus Analysis Tool', *International Journal of Corpus Linguistics* 17(3): 380–409.
- Johnson, J. and Partington, A. (2017) 'Corpus-Assisted Discourse Study of Representations of the 'Under-Class' in the English-Language Press', in E. Friginal (ed) *Studies in Corpus-Based Sociolinguistics*, London: Routledge, pp. 293–318.
- Labov, W. (1972) 'The Logic of Nonstandard English', in P. Giglioli (ed.) *Language and Social Context*, Harmondsworth: Penguin, pp. 179–215.
- Leech, G. (2002) 'Recent Grammatical Change in English: Data, Description, Theory', in K. Aijmer and B. Altenberg (eds) *Proceedings of the 2002 ICAME Conference*, Gothenburg: 61–81.
- Leech, G. (2011) 'The Modals ARE Declining. Reply to Neil Millar's 'Modal verbs in TIME: Frequency Changes 1923–2006'', *International Journal of Corpus Linguistics* 16(4): 547–64.
- Love, R., Dembry, C., Hardie, A., Brezina, V. and McEnery, T. (2017) 'The Spoken BNC2014: Designing and Building a Spoken Corpus of Everyday Conversations', *International Journal of Corpus Linguistics* 22(3): 319–44.
- MacLagan, M. A. and Hay, J. (2007) 'Getting Fed Up with Our Feet: Contrast Maintenance and the New Zealand English "Short" Front Vowel Shift', *Language Variation and Change* 19(1): 1–25.
- McEnery, T. (2005) *Swearing in English: Bad Language, Purity and Power from 1586 to the Present*, London: Routledge.
- Millar, N. (2009) 'Modal verbs in TIME', *International Journal of Corpus Linguistics* 14(2): 191–220.
- Rayson, P., Leech, G., and Hodges, M. (1997) 'Social Differentiation in the Use of English Vocabulary: Some Analyses of the Conversational Component of the British National Corpus', *International Journal of Corpus Linguistics* 2(1): 133–52.
- Schmid, H. J. (2003) 'Do Men and Women Really Live in Different Cultures? Evidence from the BNC', in A. Wilson, R. Rayson and T. McEnery (eds) *Corpus Linguistics by the Lune. Łódź Studies in Language* 8, Frankfurt: Peter Lang, pp. 185–221.
- Stubbs, M. (1996) *Texts and Corpus Analysis*, London: Blackwell.

- Subtirelu, N. (2017) 'Exploring the Intersection of Gender and Race in Evaluations of Mathematics Instructors on *Ratemyprofessors.com*', in E. Friginal (ed.) *Studies in Corpus-Based Sociolinguistics*, London: Routledge, pp. 219–35.
- Taylor, C. (2017) 'Women are Bitchy but Men are Sarcastic? Investigating Gender and Sarcasm', *Gender and Language* 11(3): 415–45.
- Torgersen, E., Kerswill, P. and Fox, S. (2006) 'Ethnicity as a Source of Changes in the London Vowel System', in F. Hinskens (ed.) *Language Variation – European Perspectives. Selected Papers from the Third International Conference on Language Variation in Europe (ICLaVE3)*, Amsterdam, June 2005, Amsterdam: Benjamins: 249–63.
- Wardhaugh, R. (2005) *An Introduction to Sociolinguistics*, London: Blackwell.