**MILLAR, NEIL. (2009)** *Modal verbs in TIME: Frequency Changes 1923-2006.* **International Journal of Corpus Linguistics.14 (2009)**

**Introduction: Why this paper?**
I was drawn to this research paper for a couple of reasons: My research interest is in English language change over time, with particular focus on American journalism (how reporting and writing styles have changed over the last 50-75 years in response to social, political, economic, and historical context). I also wanted to learn more about how to use one of the BYU corpora as a serious research tool. Although Millar's study does not address the first topic specifically (change in language with respect to American journalism per se), his use of the TIME magazine corpus caught my attention, and I also thought his paper might provide some fodder for future research ideas.

**Millar's research question—explore changes in modal verb frequency in American English**
As insight into the motivation for his study, Neil Millar (Lancaster University) discusses a diachronic study by Leech (2002, 2003), the results of which show a decline in English modal verb use. According to Millar, Leech's study focused on both British and American English using the Brown, LOB[1], FLOB[2], and Frown[3] corpora to compare changes that occur between 1961 (Brown/LOB) and 1991 (Frown/FLOB).

FLOB and Frown are both similar to Brown in that they contain "500 texts of around 2000 words each, distributed across 15 text categories, 9 informative and 6 imaginative[4]" LOB had been modeled also on Brown, but with a focus on British English from the same period (1961). (I realize you know all this, and we've come across LOB and Brown in some other papers this semester, but FLOB and Frown were new to me, so this additional context is mostly for my own edification—Millar covered the details from a high-level.)

According to Millar, Leech hypothesizes that the decline in modal verb use may be due to increasing usage of "semi-modals," but then this notion is disputed due to relatively low frequencies in the semi-modal counts (as compared to the modal verb counts—for example, although use of semi-modal "NEED to" increased by 123.2% from 1961 to 1991 in the American English data, the raw numbers are a mere 69 and 154 tokens, respectively). To say that this is somehow offsetting the modal verb "must" in the same data (1,018 tokens in 1961 and 668 tokens in 1991), even from a completely naive perspective doesn't really seem all that convincing. And so apparently, Leech ends up concluding that the decline in modal verb use might be better attributed to three factors:

---

[1] Lancaster-Oslo-Bergen Corpus of British English (1961)
[2] Freiburg-LOB Corpus of British English (1991)
[3] Freiburg update of the Brown Corpus (1991)
[4] Department of English, University of Helsinki
http://www.helsinki.fi/varieng/CoRD/corpora/FROWN/index.html

1) Americanization—tendencies for the varieties of English to shift toward American English (not really relevant for a study focused just on American English)
2) Colloquialization—trend of written norms shifting to those of speech
3) Democratization—"shift in the modal meanings that people choose to express." This notion has to do with semantic shift in the way modality is expressed. It suggests changes in the *types* of modality that society chooses to express."

All this is background to Millar's study, in which he travels down the same basic research road as Leech (with the exception, of course, that he's not looking at British English). Millar's goal in his study is three-fold:

1) replicate Leech's study insofar as examining changes in modal verb frequency;
2) compare "the findings of [Leech's] *language-in-its-entirety* approach" with his own study's "*single-genre* approach"
3) make use of "a new and still relatively unexplored linguistic resource," the TIME magazine corpus available through Brigham Young University.

**A diachronic study using the TIME Magazine Corpus**
Millar used the 'TIME Magazine Corpus (Davies 2007a) available through an online interface at http://corpus.byu.edu/." The corpus includes all 83 years of TIME magazine, starting from March, 1923 through 2006. As Millar points out, the TIME corpus differs significantly from the Brown corpora, not only in total word counts but type of texts as well. Here's a summary of the details:

| | FROWN | Brown Corpus | TIME Magazine Corpus |
|---|---|---|---|
| **Total words** | 1 million | 1 million | 100 million |
| **Text description** | 15 different categories[5] | 15 different categories | 1 single category 275,000 articles |
| **Time period** | 1991 | 1961 | 1923-2006 |
| | 1 year | 1 year | ~83 years |
| **Text words/year** | 1 million | 1 million | 1-1.5 million words |

Millar notes that the TIME Corpus has been marked-up using an automatic tagger, (CLAWS[6]), currently specified on the BYU corpora site as CLAWS7. Millar "ran queries for the same modal verbs as in the study by Leech" which were those that used syntactic criteria, specifically, modals that passed the NICE test (Huddleston and Pullum, 2002):

| Negation with –n't | Inversion with the subject | Code[7] words | Emphatic affirmation |
|---|---|---|---|
| "I can't go." | "Must you leave now?" | "You can speak Spanish, can't you?" | You will pick up your toys!" |

[5] Four sub-corpora: Press, General Prose, Learned Writing, and Fiction Informative. Press includes Reportage, Editorial, and Reviews (44, 27, and 17 texts, respectively).
[6] Constituent Likelihood Automatic Word-tagging System
[7] Code words describe as "using elliptically, as a word or code, as in tag questions."

So "contractions, including negative contractions, were included in the counts. 'll and won't treated as forms of will." Millar shows the results as standardized frequencies[8] (words per million) of modal verbs in each decade, and he calculates a straight percentage difference between the 2000s frequency as compared to the 1920s. His results are very different, in general, from Leech's.

**How Millar analyzed the data from the corpus**

Millar exported the search results from the TIME corpus into SPSS and analyzed using a "binary logistic regression model"—he plotted the data points for each year and then used the logistical model to plot the curve regression estimation (Agresti[9]). He also provides the coefficient of determination ($R^2$) values for the regression model, citing these numbers as a "goodness-of-fit test" and highlighting the fact that the $R^2$ values for most of the data are well above a 50% threshold. These are the $R^2$ values for the five modal verbs that show the greatest decline (or increase) over the time period from 1923 to 2006.

| modal verb | shall | can | may | must | could |
|---|---|---|---|---|---|
| **$R^2$** | 0.924 | 0.819 | 0.675 | 0.665 | 0.612 |

Miller's results show a decline in the frequency of *shall* and *must*, and an increase in frequency of modal verbs *can*, *could*, and *may*. Although Leech's results also showed decline in *shall* and *must*, the results for *can, could,* and *may* are considerably different. For example, where Leech recorded a 32.4% decrease in the use of *may* in the Brown Corpus, Millar finds a 59.7% increase.

Taking into account corpus type and timeframe differences, Millar conducts some additional analysis, comparing just the *Press* sub-corpora of Frown to Brown, and also comparing just the years 1961 to 1991 of the TIME corpora. Here's a recap of some of the results from these additional comparisons.

| | Brown-Frown 1961-1991 | | TIME Corpus | |
|---|---|---|---|---|
| | **[1] All** | **[2] Press** | **[3] 1920s-2000s** | **[4] 1961-1991** |
| shall | -43.8% | -70.2% | -95.6% | -87.3% |
| can | -1.5% | +1.1% | +113.4% | -3.1% |
| may | -32.4% | -39.3% | +59.7% | +54.0% |
| must | -34.4% | -33.9% | -48.4% | -34.9% |
| could | -6.8% | +25.9% | +24.0% | +103.0% |

As Millar points out, restricting the analysis to just the Press sub-corpora seems to bring the results closer together (for example, *could*) but the modal verb *may* shows an even

---

[8] Millar does not provide any information about his approach to the "standardized values," or if he simply picked up the values provided by the corpus, assuming he selected that "per million" as one of the filters on his queries.

[9] Agresti, A. 2002. Categorical Data Analysis, 2nd Edition. New York. Wiley-Interscience.

greater decline than in the corpus as a whole. Based on this, Millar notes that "text type …is not wholly satisfactory as an explanation for the observed discrepancies" [between his data and Leech's], and he concludes that the methodology by which these figures of percentage change have been obtainer may be flawed…. the simple comparison of two data points from TIME …produces a picture of change that contradicts the overall pattern…"

Beyond the frequency comparisons, Millar conducted a "selected analysis of some semantic aspects of modal verbs," taking random samples for modal verbs *may, should,* and *must* from 1923, 1960, and 2000 and semantically coding a "minimum of one-third of all concordance lines" for each of these years. None of the implementation details are provided, but he does find that the sense of the words (in terms of the modality being expressed) is changing. For example, the word *may* 'shows a substantial shift towards epistemic sense of possibility, the expression of the writer's 'lack of confidence in the truth of the proposition' (Coates 1983: 133)"—the trend is line with Leech's findings. In 1923, "56% of all instances of *may*…use the epistemic sense" and this rises "to 78% by 1960 and 94% by 2000" according to Millar's analysis. His data shows changes in senses over time for most all the modal verbs.

**Conclusions and Implications of Millar's study**
Millar concurs with Leech, in that "changes in modal verbs are taking place," and that "changes are most evident for *shall, must,* and *ought* (decrease) and *may, can,* and *could* (increase)." He suggests that "notions of colloquialization, democratization and stylistic change go some way to explain these changes," but he also suggests that there could be "an increase in speculation on the part of journalists, resulting in greater use of markers of tentative epistemic modality…"

The implications he raises have to do with methodological approaches, and whether it's sufficient to extrapolate from two data points. The two corpora (Brown/Frown) and TIME seem so different in terms of content that it seems wrong-headed to even try to compare data from one to the other and Millar cautions against drawing any conclusions, noting "the extension of these generalizations to the English language as a whole will necessarily remain tentative."

On the other hand, he points out, an "advantage of the single-source approach taken in this study (*all* articles published… over an 83 year period) is that these generalizations are based on data which are robust." Further, mentioned early on in the paper is that TIME is the most-read newsmagazine in America (Hau 2008), with a weekly circulation of c. four million copies in 2006, generalizations from the TIME corpus to the genre of newsmagazine can be made with reasonable confidence."

However, I ended up being a bit disappointed with this paper—it seemed to try to cover too much ground, with very little coverage given to some of the low-level details of how he actually did the analysis. I got very sidetracked trying to duplicate some of his results, unsuccessfully. In the meantime, I've found what looks to be a good paper by Mark Davies ("*The 385+ million word Corpus of Contemporary American English—Design, architecture, and linguistic insights" (2009)*), so I'll leave Millar here and move on to Davies.