Check for updates

# Phonological processes of consonants from orthographic to pronounced words in the Buckeye Corpus*

Byunggon Yang**

*Department of English Education, Pusan National University, Pusan, Korea*

## Abstract

This paper investigates the phonological processes of consonants in pronounced words in the Buckeye Corpus and compares the frequency distribution of these processes to provide a clearer understanding of conversational English for linguists and teachers. Both orthographic and pronounced words were extracted from the transcribed label scripts of the Buckeye Corpus. Next, the phonological processes of consonants in the orthographic and pronounced labels were tabulated separately by onsets and codas, and a frequency distribution by consonant process types was examined. The results showed that the majority of the onset clusters were pronounced as the same sounds in the Buckeye Corpus. The participants in the corpus were presumed to speak semiformally. In addition, the onsets have fewer deletions than the codas, which might be related to the information weight of the syllable components. Moreover, there is a significant association and strong positive correlation between the phonological processes of the onsets and codas in men and women. This paper concludes that an analysis of phonological processes in spontaneous speech corpora can contribute to a practical understanding of spoken English. Further studies comparing the current phonological process data with those of other languages would be desirable to establish universal patterns in phonological processes.

Keywords: phonological processes, consonants, sex, orthographic, pronounced, Buckeye Corpus

## 1. Introduction

A spontaneous speech corpus stores much information on how Americans communicate in their daily lives. The phonological processes in the corpus are useful for linguists and teachers to examine and propose practical guidelines for fluent English pronunciation. However, not many studies have reported on the authentic processes revealed in the pronunciation corpus, probably because of the tremendous amount of information in the corpus to derive general trends from the processes. Any manual processes of classifying the data into different categories would take a long time, let alone the errors and mistakes arising from a tedious task. Fortunately, recent software is available to handle the main part of the task. Careful programming can automate the whole procedure, but it may demand much more time and effort to deal with issues arising from complex phonological processes and several redundant usages of human speech. Additional manual work may save some time to resolve these issues.

In English, the medial phoneme /t/ is reported to be realized as three forms: its canonic form with a closure and release; a glottal stop; and a flap (Patterson & Connine, 2001). From an analysis of telephone conversation corpora, Patterson & Connine (2001) reported that 96.4% of the tokens with a medial /t/ for the North Midland dialect were pronounced as a flapped /t/ followed by a glottal stop for 2.12% and the canonic /t/ for 1.48%. The flapping occurred more in the intimate conversation. Ernestus, Hanique and Verboom (2015) examined the effect of the formality of the speech situation on the degree of reduced pronunciation variants. They reported that the read-aloud stories showed the lowest percentage of semantically weak variants, but a higher percentage emerged in casual face-to-face and telephone conversations.

A syllable consists of an onset, nucleus, and coda. Yang (2016) reviewed English syllable components and previous syllabification procedures and explored the phoneme distribution and syllable structure of entry words in the Carnegie Mellon University English Pronouncing Dictionary. He found that English words contained more consonants than vowels, with a 6:4 frequency ratio of consonants to vowels. He also noted that the frequency distribution of consonants was comparable to that in the Buckeye Corpus. Yang (2012) examined 2.6 million consonants and vowels of orthographic and pronounced words in the Buckeye Corpus and reported that the orthographic symbols were massively reduced in the pronounced words in the conversation. The rate of reduction amounted to 38.2%. Stops, fricatives, and nasals accounted for 75% of the consonant inventory. His study pooled the consonants without considering the syllable components. As Kessler & Treiman (1997) noted the asymmetry of onsets and codas in Spanish, the syllable components would better represent the distribution of consonants in a language. They indicated that Spanish had many codas with the anterior coronals /θ, ð, s, l, n, r/ but not any other consonants.

On the other hand, Yang's (2018) recent study on the Buckeye Corpus investigated the phonological processes of English vowels and found that the speakers tended to keep their vowel quantity in spontaneous speech. He reported that 95% of the orthographic vowels were pronounced as the same number of vowels, while the remaining 5% were inserted or deleted. He also suggested that the speakers preferred deletion to insertion in the phonological processes of vowels. He attributed the phenomenon to the effectiveness of deletion over insertion in the reduction of efforts for their production and speech rates. Moreover, the vowel quality varied over the level of chance: 58.2% of the vowels in orthographic words were pronounced with the same vowels, whereas 41.8% of them were pronounced with different vowels. The vowels for communication, i.e., stressed syllable and one-syllable content words, were mostly preserved, while the vowels in unstressed and function words were under phonological processes. This study was motivated by the dearth of papers on the phonological processes of English consonants in a large spontaneous speech corpus comparable to the study of English vowels reviewed above.

The main purpose of this study was to examine the phonological processes of English consonants on the basis of syllable components in a spontaneous speech corpus of American English. Specifically, the current study was designed to investigate distributions of the phonological processes of English consonants in the Buckeye Corpus by onsets and codas and to compare sex-related differences in these processes.

## 2. Method

### 2.1. The Buckeye Corpus

The Buckeye Corpus consisted of forty speakers (see Pitt et al., 2007 for details). They were divided into two sex groups, 20 men and 20 women, and two age groups, a younger group under 30 and an older group over 40. All participants were speakers of American English born in Columbus, Ohio. They passed a dialect screening test through a telephone interview. The recruited participants joined one-hour-long individual sessions with a conversation on general topics, including personal backgrounds and views on miscellaneous issues.

### 2.2. Data Collection and Analysis

The phonetically transcribed text files of the Buckeye Corpus were divided into sex and age group folders. Then, an R (R Core Team, 2019) script was created to collect an integrated text file to which all the files within each folder were appended as character strings (see Yang, 2018 for details). The author removed format errors and unrelated labels from the integrated file in Microsoft Excel. The total number of words in the file was 283,522. The number of consonants for the orthographic transcriptions alone amounts to 1,580,547 in 30 phonetic categories.

The integrated Excel file was imported onto the R Studio, and another R script was created to trace the phonological processes from the orthographic to the pronounced phonetic symbols of English consonants. The procedure was based on the syllabification rules prescribed in Noyer (2016) and the list of English onset clusters in Duanmu (2002) and Williamson (2014). In brief, the script initialized the Buckeye list of 33 vowels and 27 single consonants along with 54 two-symbol onset clusters and 9 three-symbol onset clusters. Then, a template matrix was created with one row and forty columns. Phonetic strings of the orthographic and pronounced words were extracted along with the number and position of vowels in the words. The syllabification rules were applied to find the onsets and codas of the words. When the number of vowels in a given orthographic word matched that of the pronounced one, one of the columns in the matrix was recorded as "same"; otherwise, it was recorded as "different", which indicates variants. Some variants were syllabic deletions, such as the word "probably". Its orthographic phonetic form was transcribed as [p r aa b ah b l iy], while one of the pronounced ones was noted as [p r aa l iy], in which the stressed first syllable was pronounced but the second syllable and one of the onset cluster of the third syllable were deleted. The different rows were manually checked to assign appropriate consonantal processes according to the corresponding components in the given syllable structure later. The number of vowels in a word was used to set a looping frequency in the script. Generally, the syllabification rules were applied to obtain the onset (cluster) first, and then the remaining consonants between adjacent syllables were assigned to the current or previous syllable as the coda. The same procedure was applied to the pronounced words, and the list for each syllable was appended to the result file. In addition, 14,231 (5% of all words) different rows of the file were manually checked to match the phonological processes from the orthographic to the corresponding pronounced syllables appropriately.

A third R script was created to combine rows of onsets or codas

**Table 1.** Phonological processes of single onsets in the Buckeye Corpus. The column name *orth* represents the orthographic forms, while *pron* stands for the pronounced phonetic symbols of English consonants. *n* is the number of occurrences, and ∅ represents deletion.

| Orth | Pron | n | Orth | Pron | n | Orth | Pron | n | Orth | Pron | n | Orth | Pron | n | Orth | Pron | n |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| b | b | 13,507 | dh | dh | 19,281 | jh | jh | 3,585 | n | n | 10,933 | t | t | 11,316 | w | w | 19,670 |
|  | ∅ | 1,026 |  | n | 3,295 |  | ch | 164 |  | nx | 4,069 |  | dx | 4,052 |  | ∅ | 1,468 |
|  | v | 145 |  | ∅ | 3,246 |  | zh | 163 |  | ∅ | 1,251 |  | ∅ | 1,655 |  | hhw | 125 |
|  | m | 92 |  | th | 996 |  | sh | 79 | p | p | 7,380 |  | d | 1,442 |  | mw | 61 |
|  | p | 69 |  | d | 577 |  | ∅ | 77 |  | b | 118 |  | nx | 585 | y | y | 12,435 |
| ch | ch | 1,620 |  | z | 364 |  | d | 55 | r | r | 7,861 |  | tq | 420 |  | ∅ | 857 |
|  | sh | 366 |  | t | 303 | k | k | 11,264 |  | ∅ | 1,701 |  | n | 378 |  | ny | 170 |
| d | d | 9,383 |  | s | 231 |  | kl | 145 | s | s | 14,251 |  | ch | 159 |  | sh | 122 |
|  | dx | 2,907 |  | l | 184 |  | g | 139 |  | z | 80 |  | r | 72 |  | zh | 115 |
|  | ∅ | 1,157 |  | dx | 139 |  | ∅ | 99 |  | sh | 70 |  | tr | 72 |  | ch | 88 |
|  | n | 246 | g | g | 7,321 | l | l | 14,060 |  | ∅ | 60 |  | s | 50 |  | jh | 61 |
|  | nx | 138 |  | k | 99 |  | ∅ | 1,320 | sh | sh | 4,065 | th | th | 4,716 |  | n | 53 |
|  | t | 109 |  | ∅ | 68 |  | kl | 89 |  | ch | 66 |  | dh | 279 | z | z | 1,712 |
|  | tq | 80 | hh | hh | 11,357 | m | m | 12,919 | v | v | 4,581 |  | ∅ | 174 |  | s | 168 |
| f | f | 6,807 |  | ∅ | 2,751 |  | ∅ | 190 |  | ∅ | 181 |  | t | 67 |  | ∅ | 68 |
|  | fr | 321 |  |  |  |  |  |  |  | f | 60 |  |  |  | zh | zh | 177 |

using the R function "rbind" for all speakers or for the groups of men and women and to write the table data counting the number of phonological processes. In addition, chi-squared tests and correlation analyses were conducted to determine the association of the frequency distribution of the phonological processes of consonants by sex and the degree of correlation between the two groups.

## 3. Results and Discussion

### 3.1. Phonological Processes of Onsets

Table 1 lists the frequency distribution of the phonological processes of single onsets in the Buckeye Corpus. The list includes only cases with more than 50 variants to make the table simpler. In addition, consonant insertions were excluded. The same screening criteria were applied to create the following tables in this paper. Consonant insertion almost never occurred compared with the abundant deletions. For example, there were 3,246 /dh/ deletions and 2,751 /hh/ deletions for the onsets. On the other hand, there were 131 /l/ insertions followed by 114 /n/, 61 /r/ and 61 /w/ deletions for the onsets, while there were 822 /r/ insertions followed by 809 /l/, 152 /n/, 110 /s/, and 73 /ng/ insertions for the codas.

Generally, the table shows that the majority of the onset consonants were pronounced as the same orthographic symbols and with the same voicing quality in the conversation. Specifically, 19,670 /w/s were pronounced as the same consonants, followed by 19,291 /dh/s and 14,251 /s/s and 14,060 /l/s. The speakers in the Buckeye Corpus must have pursued this as a strategy to ensure that they would be clearly understood. We will discuss the important role of onsets for word identification in the following section with additional coda data.

The alveolar stops /d, t/ have relatively diverse variants. For example, /t/ lists 11 variants followed by the voiced fricative /dh/ and its voiced counterpart /d/. The palatalized /ch/ occurs only 159 instances, while its voiced counterpart /jh/ occurs 23 instances. These results might be related to the semiformal mode of the interview, which we will discuss below. On the other hand, the voiced stops in the orthographic forms tend to keep the voice quality except for the dental stop /t/; there are 1,442 instances of voiced /d/. The fricatives also seem to maintain the same sound quality and voicing quality. The occurrence of 996 instances of /dh/ to /th/ might be related to the high occurrence of the fricative itself. For example, the proportion of the "same" consonants in both the orthographic and pronounced forms (dh<dh) to the different consonants (dh<th) was approximately 5%. The voiceless counterpart /th/ accounts for approximately 6%. This paper uses the percentage of the number of the same consonant processes as that of the given variants to estimate the relative weight of variants. The deletion of the fricative /dh/ topped the list, followed by the fricative /hh/. The missing fricatives are mostly in the middle unstressed syllables, as can be seen in vowel studies (Yang, 2018). The deletion of the /l/ sound also accounts for approximately 10% of the same pronounced forms.

It is interesting to see the flapped variants /dx/ of /t/, which amount to 4,052 instances along with 585 occurrences of /nx/. Because the Buckeye Corpus has relatively fewer reductions or variants, the speakers were presumed to talk semiformally. Zue & Laferriere (1979) listed six phonetic environments of the phoneme /t/ in the medial position of English words. More detailed analyses of the variants that consider the phonetic environments are desirable to complete the study of these phonological processes. Nonnative speakers such as Koreans may need to pay attention to those variants in pronouncing and listening exercises in order to carry out a casual conversation with Americans.

The nasal /m/ lists the same process and deletion, while /n/ lists an additional nasal flap /nx/. Interestingly, no /m/ variants occurring over 50 instances were reported here. On the other hand, the rate of the same consonant to the nasal flap of the orthographic /n/ is approximately 37%, which seems remarkable. In addition, 1,251 /n/ deletions occur in the spontaneous speech. The reason for the robust use of /m/ might be related to the distribution of English consonants in the place of articulation: English has fewer labial consonants and variants available than alveolar consonants and variants. Here, the labial /m/ occurs more in the onset than the alveolar /n/ does. The labial /w/ lists a few variants, and its deletion rate amounts to approximately 7.5%. The deletion rate of /r/ is 21.6%, while that for /l/ is 9.4%. We may say that the lateral /l/ is more robust to changes in the phonological process. The glide /y/ has many variants, which might be related to the fact that the glide tends to be combined with the preceding codas or the following vowels to form consonant clusters or diphthongs. On the other hand, we also note that the voiceless fricatives /f, s, sh/ prevail in the onset positions. Among them, the fricative /s/ is the most prominent. We will return to the frequency distribution of consonants by syllable components in the

following coda section.

| Orth | Pron | n | Orth | Pron | n | Orth | Pron | n |
|------|------|------|------|------|------|------|------|------|
| bl | bl | 597 | hhy | hhy | 51 | sw | sw | 86 |
| bl | ∅ | 151 | ly | ly | 70 | thr | thr | 543 |
| bl | l | 104 | ly | y | 50 | thr | th | 53 |
| bl | b | 103 | my | my | 154 | tr | tr | 476 |
| br | br | 470 | ny | ny | 218 | tr | chr | 425 |
| by | by | 68 | pl | pl | 758 | tr | r | 131 |
| dr | dr | 382 | pr | pr | 1,535 | tw | tw | 294 |
| dr | jhr | 102 | pr | p | 175 | vr | vr | 570 |
| dy | jh | 77 | py | py | 152 | vr | r | 126 |
| fl | fl | 217 | sf | sf | 57 | vr | v | 60 |
| fr | fr | 987 | sk | sk | 1,022 | vy | vy | 59 |
| fr | f | 191 | sl | sl | 236 | skr | skr | 73 |
| fy | fy | 113 | sm | sm | 196 | spl | spl | 57 |
| gl | gl | 76 | sp | sp | 847 | str | str | 234 |
| gr | gr | 979 | st | st | 2,873 | str | schr | 92 |
| kl | kl | 760 | st | s | 133 | str | r | 91 |
| kr | kr | 513 | st | ch | 67 | | | |
| kw | kw | 321 | | | | | | |
| ky | ky | 146 | | | | | | |

Table 2 lists the phonological processes of the onset clusters. The majority of the onset clusters have fewer variants than the single onsets and almost no clusters with more than 50 deletions except the cluster /bl/. That cluster lists four variants and 151 deletions. The cluster /tr/ has almost comparable variants of /tr/ and /chr/. The voiced counterpart /dr/ changed to /jhr/. The variants are palatalized alveolars. The fewer variants may account for the strong association of the two or three single consonants. If one of the onset clusters changes, listeners may have difficulty guessing the words. The cluster /pr/ lists 1,535 counts, while /pl/ lists 758 counts. A similar trend of a preferred combination of a consonant plus /r/ can be seen between /fr/ (987) and /fl/ (217) and between /gr/ (979) and /gl/ (76). However, the cluster /br/ lists 470 instances against 597 instances of /bl/. It seems very difficult to generalize the trend because the occurrence of clusters largely depends on both the phonological feature of the clusters themselves and the frequency of everyday expressions with the clusters. On the other hand, we note that the distribution of the onset triple clusters in the inventory varies. The cluster /str/ is most prevalent, with 234 instances, followed by /skr/; the 47 instances of /spr/ are not listed on the table. Whether the uneven distribution is related to the perceptual prominence of the clusters would be interesting to study in the future.

| Orth | Pron | n | Orth | Pron | n | Orth | Pron | n | Orth | Pron | n |
|------|------|------|------|------|------|------|------|------|------|------|------|
| b | b | 352 | l | l | 5,897 | p | p | 1,893 | th | th | 1,251 |
| b | ∅ | 82 | l | ∅ | 1,892 | p | ∅ | 71 | th | dh | 344 |
| ch | ch | 1,331 | m | m | 10,864 | r | r | 8,117 | th | ∅ | 87 |
| d | d | 6,402 | m | em | 848 | r | ∅ | 5,527 | v | v | 5,687 |
| d | dx | 1,474 | m | ∅ | 364 | s | s | 4,000 | v | ∅ | 2,362 |
| d | ∅ | 819 | m | n | 66 | s | sh | 148 | v | f | 622 |
| d | b | 94 | n | n | 16,833 | s | z | 133 | v | b | 96 |
| d | tq | 67 | n | ∅ | 4,237 | s | ∅ | 69 | v | z | 11,267 |
| d | t | 57 | n | nx | 1,961 | sh | sh | 285 | z | s | 2,280 |
| f | f | 2,713 | n | m | 327 | t | tq | 9,202 | z | zh | 252 |
| f | v | 64 | n | ng | 160 | t | t | 9,116 | z | ∅ | 220 |
| f | ∅ | 63 | n | nch | 72 | t | dx | 4,954 | z | sh | 96 |
| g | g | 629 | ng | ng | 5,884 | t | ∅ | 3,844 | z | rz | 50 |
| g | ∅ | 68 | ng | n | 1,333 | t | d | 718 | | | |
| jh | jh | 408 | ng | ∅ | 906 | t | n | 174 | | | |
| k | k | 8,559 | ng | nx | 254 | t | p | 95 | | | |
| k | ∅ | 258 | ng | em | 85 | t | k | 78 | | | |
| k | g | 168 | ng | eng | 78 | t | b | 75 | | | |
| | | | ng | m | 71 | | | | | | |

## 3.2. Phonological Processes of Codas

Table 3 presents the phonological processes of single codas. The deletion rate of the lateral /l/ amounts to 32.1% of the same consonant processes. Similarly, we can see large deletions of /n, v/. Moreover, the number of variants in the codas is smaller than the number in the onsets. The alveolar consonant /t/ lists nine variants followed by seven /ng/ variants. This result might be related to the inclusion of all the onsets which occur in the initial, middle, and final syllables. Moreover, the number of deletions in the onsets is relatively smaller than that of the codas. For example, 14.6% of the onset /t/ to the same consonant process was deleted, while 42.2% of the coda /t/ was deleted. This trend may be attributed to the information weight of the syllable components. Generally, the onsets seem more important in the delivery of ideas than the codas. Cutler (1982) described the beginning of a word as the most salient part for identification. In other words, speakers immediately guess words without waiting for their final coda. In addition, Yang (2018) mentioned that speakers tended to preserve vowels that were important to convey their thoughts but tended to change or delete vowels that were relatively unimportant to deliver their intended message. The same generalization may be applicable to the current results.

We described the phonological processes, but here let us briefly examine the consonant distribution in light of voicing and syllable components. Table 3 shows that the voiced /z/ prevails over the voiceless /s/, which occurs almost three times. The same trend can be seen between the fricatives /f, v/. It is noticeable that the onsets in Table 1 show the reverse trend: the voiceless consonants prevail. On the other hand, the voiceless stops in the coda still prevail. The stop /p/ occurs 5.4 times that of its voiced counterpart /b/. It seems difficult to generalize this conflicting distribution among the onsets and codas in the current corpus. Further analysis of more spontaneous speech corpora may be necessary to draw a conclusive statement on the phenomenon. We can also note that the voiceless stops /p, t, k/ occur much more frequently than their voiced counterparts. The voiced /b/ and /g/ occur 352 and 629 instances, respectively. Maddieson (1984) attempted to establish universals

across world languages and reported that 80% or more languages had the consonants /p, t, k, m, n, s, j/. In addition, Eckman (2004) reviewed the Markedness Hypothesis and posited that the notion of typological markedness might be closely related to language universals (see the review in Yang, 2012). Voiceless consonants occur in languages more often than voiced consonants. Thus, voiceless consonants are considered simpler and more natural than voiced ones. However, we note that the distribution of consonants is different in onsets and codas. Thus, an observation of the total inventory distribution of consonants may have to be reexamined in light of the syllable components. In other words, if we combine the occurrence of the onsets and codas into one consonant category, the frequency of voiced and voiceless consonants may offset and obscure the individual characteristics of the given consonant category according to the syllable components.

Table 4 lists the phonological processes of coda clusters in the corpus. The coda clusters have many more variants than the onset clusters have. Additionally, relatively greater deletions occur in the clusters /nd, nt, st, kt/. The rate of processes from the clusters /nd/ and /nt/ to a single coda /n/ amounts to 26.9%. The first component of the cluster /st/ remains the same in 2,353 instances, while 1,621 /st/ instances remain the same in the pronounced forms. Many coda clusters keep their orthographic forms in the pronounced forms. Homorganic consonant clusters are expected to be deleted in the pronounced forms because it would be rather difficult for the participants to produce the clusters with a different manner of production almost simultaneously. For example, the majority of /nd, nt/ and /st/ become one /n/ and /s/, respectively. However, the similar homorganic clusters /ld/, /lt/, /lz/, /ls/, and /ts/ tend to remain the same.

Here, we discuss the frequency distribution of the clusters. Nasal, alveolar and lateral flaps are prevalent in the coda clusters. The majority of the clusters start with the consonant /n/ followed by /l/ or /r/. Among them, the consonant /n/ is most favored as the first component of the clusters. The frequency distribution of English consonants may have to consider these consonant clusters without separating them into single consonants. The frequency distribution of the cluster /ts/ is greater than that of /st/. In Table 2, we observed 2,873 onset /st/ clusters. Here, the mirror image of /ts/ has 3,346 instances. Clusters consisting of the fricative /s/ followed by other consonants are considered problematic for sonority analysis. The sonority sequencing principle states that "only sounds of higher sonority rank are permitted between any member of a syllable and the syllable peak (Clements, 1990:285)". However, we have seen that the clusters /st/ and /ts/ do not follow the symmetry in syllable structure. In other words, the fricative /s/ is considered more prominent than the stop /t/ in a sequence. The statistical distribution of /st/ in the corpus might provide a gradient explanation of the asymmetry depending on the syllable components, which will be examined in the future.

**Table 4.** Phonological processes of coda clusters in the Buckeye Corpus. The column name *orth* represents the orthographic forms, while *pron* stands for the pronounced phonetic symbols of English consonants. *n* is the number of occurrences, and ∅ represents deletion.

| Orth | Pron | n | Orth | Pron | n | Orth | Pron | n | Orth | Pron | n |
|---|---|---|---|---|---|---|---|---|---|---|---|
| dz | dz | 431 | nd | n | 7,250 | ngk | ngk | 1,417 | rt | rdx | 226 |
|  | ds | 277 |  | ∅ | 2,936 |  | ng | 234 |  | rt | 195 |
| ft | ft | 84 |  | nd | 1,429 |  | k | 88 |  | rtq | 192 |
| ks | ks | 812 |  | nx | 1,198 | ngz | ngz | 583 |  | r | 97 |
| kt | kt | 495 |  | m | 133 |  | ngs | 175 | rts | rts | 104 |
|  | k | 204 |  | nt | 100 | njh | njh | 104 | rz | rz | 902 |
|  | ∅ | 107 |  | d | 70 | nth | nth | 77 |  | rs | 203 |
| kts | ks | 55 |  | em | 52 | nts | nts | 233 |  | z | 170 |
| ld | ld | 478 | ns | ns | 436 |  | ns | 113 | sk | sk | 62 |
|  | l | 125 |  | s | 73 |  | ts | 89 | st | s | 2,353 |
|  | ldx | 104 |  | nts | 57 |  | s | 58 |  | st | 1,621 |
| lf | lf | 173 | nt | d | 73 | nz | nz | 600 |  | ∅ | 79 |
| lp | lp | 140 |  | n | 1,952 |  | ns | 174 |  | sh | 65 |
| ls | ls | 160 |  | ∅ | 927 |  | z | 105 |  | z | 54 |
| lt | lt | 56 |  | tq | 920 | ps | ps | 131 | ts | ts | 3,346 |
| lz | lz | 271 |  | nt | 772 | pt | pt | 90 |  | s | 1,303 |
|  | ls | 74 |  | nx | 464 |  | p | 68 |  | tqs | 72 |
| mz | mz | 494 |  | ntq | 396 | rd | rd | 320 |  | z | 70 |
|  | ms | 158 |  | t | 232 |  | rdx | 59 | zd | s | 124 |
| nch | nch | 86 |  | m | 72 | rk | rk | 68 |  | z | 116 |
|  |  |  | ndz | nz | 146 | rs | rs | 140 |  | zd | 111 |
|  |  |  |  | ndz | 102 |  |  |  |  |  |  |

### 3.3. Comparison of the Phonological Processes by Sex

Tables 5 and 6 list the frequency counts of the phonological processes of men and women in the Buckeye Corpus. Instead of listing all the processes exhaustively, 94 instances of the onset and coda processes were chosen to conduct a statistical analysis of the association between the two sex groups. The usual approach would be to conduct the test over a large table of all phonological processes. However, some slots would have too few occurrences for the standard chi-squared test.

**Table 5.** Major frequency distribution of the onsets of men and women in the Buckeye Corpus. The column head *orth* represents orthographic forms, while *pron* stands for the pronounced phonetic symbols of English consonants. *n* is the number of occurrences. *m_n* and *w_n* are the numbers of occurrences of men and women, respectively. ∅ represents deletion.

| Orth | Pron | m_n | w_n | Orth | Pron | m_n | w_n | Orth | Pron | m_n | w_n |
|---|---|---|---|---|---|---|---|---|---|---|---|
| b | b | 6,902 | 6,605 | jh | ch | 92 | 72 | st | s | 79 | 54 |
|  | v | 82 | 63 |  | jh | 1,735 | 1,850 |  | st | 1,611 | 1,262 |
|  | ∅ | 517 | 509 |  | zh | 90 | 73 | str | str | 123 | 111 |
| bl | bl | 318 | 279 | k | g | 80 | 59 | t | ch | 91 | 68 |
| br | br | 280 | 190 |  | k | 5,847 | 5,417 |  | d | 711 | 731 |
| ch | ch | 772 | 848 |  | kl | 77 | 68 |  | dx | 2,075 | 1,977 |
|  | sh | 161 | 205 | kl | kl | 409 | 351 |  | n | 190 | 188 |
| d | d | 4,814 | 4,569 | kr | kr | 274 | 239 |  | nx | 332 | 253 |
|  | dx | 1,439 | 1,468 | kw | kw | 193 | 128 |  | t | 5,932 | 5,384 |
|  | n | 128 | 118 | ky | ky | 76 | 70 |  | tq | 178 | 242 |
|  | nx | 60 | 78 | l | l | 7,443 | 6,617 |  | ∅ | 987 | 668 |
|  | ∅ | 614 | 543 |  | ∅ | 687 | 633 | th | dh | 188 | 91 |
| dh | d | 341 | 236 | m | m | 6,634 | 6,285 |  | th | 2,135 | 2,581 |
|  | dh | 9,586 | 9,695 |  | ∅ | 128 | 62 | thr | thr | 293 | 250 |
|  | dx | 79 | 60 | my | my | 97 | 57 | tr | chr | 232 | 193 |
|  | l | 129 | 55 | n | n | 5,572 | 5,361 |  | tr | 250 | 226 |
|  | n | 1,565 | 1,730 |  | nx | 1,967 | 2,102 | tw | tw | 156 | 138 |
|  | s | 118 | 113 |  | ∅ | 753 | 498 | v | v | 2,290 | 2,291 |
|  | t | 218 | 85 | ny | ny | 110 | 108 |  | ∅ | 124 | 57 |
|  | th | 483 | 513 | p | p | 4,073 | 3,307 | vr | vr | 259 | 311 |
|  | z | 194 | 170 | pl | pl | 453 | 305 | w | hhw | 53 | 72 |
| dr | dr | 192 | 190 | pr | p | 92 | 83 |  | w | 10,188 | 9,482 |
| f | f | 3,708 | 3,099 |  | pr | 855 | 680 |  | ∅ | 906 | 562 |
|  | fr | 162 | 159 | py | py | 84 | 68 | y | ny | 80 | 90 |
| fl | fl | 116 | 101 | r | r | 3,879 | 3,982 |  | sh | 53 | 69 |
| fr | f | 95 | 96 |  | ∅ | 1,015 | 686 |  | y | 5,974 | 6,461 |
|  | fr | 567 | 420 | s | s | 7,789 | 6,462 |  | ∅ | 634 | 223 |
| g | g | 3,855 | 3,466 | sh | sh | 2,024 | 2,041 | z | s | 66 | 102 |
| gr | gr | 477 | 502 | sk | sk | 519 | 503 |  | z | 932 | 780 |
| hh | hh | 5,373 | 5,984 | sl | sl | 143 | 93 | zh | zh | 86 | 91 |
|  | ∅ | 1,382 | 1,369 | sm | sm | 107 | 89 |  |  |  |  |
|  |  |  |  | sp | sp | 454 | 393 |  |  |  |  |

**Table 6.** Major frequency distribution of the codas of men and women in the Buckeye Corpus. The column head *orth* represents orthographic forms, while *pron* stands for the pronounced phonetic symbols of English consonants. *n* is the number of occurrences. *m_n* and *w_n* are the numbers of occurrences of men and women, respectively. ∅ represents deletion.

| Orth | Pron | m_n | w_n | Orth | Pron | m_n | w_n | Orth | Pron | m_n | w_n |
|---|---|---|---|---|---|---|---|---|---|---|---|
| b | b | 203 | 149 | nd | n | 3,368 | 3,882 | rdx | rdx | 125 | 101 |
| ch | ch | 672 | 659 |  | nd | 793 | 636 | rt | rt | 107 | 88 |
| d | d | 3,242 | 3,160 |  | nx | 584 | 614 |  | rtq | 78 | 114 |
|  | dx | 730 | 744 |  | ∅ | 1,691 | 1,245 | rz | rs | 78 | 125 |
|  | ∅ | 456 | 363 | ndz | ndz | 52 | 50 |  | rz | 532 | 370 |
| dz | ds | 75 | 202 |  | nz | 87 | 59 |  | z | 117 | 53 |
|  | dz | 200 | 231 | ng | n | 752 | 581 | s | s | 2,155 | 1,845 |
| f | f | 1,426 | 1,287 |  | ng | 3,043 | 2,841 |  | sh | 95 | 53 |
| g | g | 339 | 290 |  | nx | 156 | 98 |  | z | 75 | 58 |
| jh | jh | 218 | 190 |  | ∅ | 640 | 266 | sh | sh | 149 | 136 |
| k | g | 86 | 82 | ngk | ng | 97 | 137 | st | s | 1,081 | 1,272 |
|  | k | 4,723 | 3,836 |  | ngk | 545 | 872 |  | st | 804 | 817 |
|  | % | 170 | 88 | ngz | ngs | 75 | 100 | t | d | 333 | 385 |
| ks | ks | 397 | 415 |  | ngz | 282 | 301 |  | dx | 2,551 | 2,403 |
| kt | k | 118 | 86 | ns | ns | 230 | 206 |  | n | 68 | 106 |
|  | kt | 268 | 227 | nt | n | 956 | 996 |  | t | 5,206 | 3,910 |
| l | l | 3,009 | 2,888 |  | nt | 409 | 363 |  | tq | 3,907 | 5,295 |
|  | ∅ | 941 | 951 |  | ntq | 157 | 239 |  | ∅ | 2,191 | 1,653 |
| ld | ld | 252 | 226 |  | nx | 245 | 219 | th | dh | 201 | 143 |
| lf | lf | 93 | 80 |  | t | 151 | 81 |  | th | 609 | 642 |
| ls | ls | 79 | 81 |  | tq | 371 | 549 | ts | s | 670 | 633 |
| lz | lz | 138 | 133 |  | ∅ | 596 | 331 |  | ts | 1,847 | 1,499 |
| m | em | 550 | 298 | nts | ns | 61 | 52 | v | f | 289 | 333 |
|  | m | 5,257 | 5,607 |  | nts | 102 | 131 |  | v | 2,914 | 2,773 |
|  | ∅ | 253 | 111 | nz | ns | 75 | 99 |  | ∅ | 1,324 | 1,038 |
| mz | ms | 62 | 96 |  | nz | 346 | 254 | vd | vd | 72 | 99 |
|  | mz | 287 | 207 | p | p | 1,021 | 872 | z | s | 958 | 1,322 |
| n | m | 154 | 173 | ps | ps | 64 | 67 |  | z | 5,770 | 5,497 |
|  | n | 8,739 | 8,094 | r | r | 4,446 | 3,671 |  | zh | 120 | 132 |
|  | ng | 72 | 88 |  | ∅ | 2,998 | 2,529 |  | ∅ | 135 | 85 |
|  | nx | 1,111 | 850 | rd | rd | 181 | 139 |  |  |  |  |
|  | ∅ | 2,541 | 1,696 | rs | rs | 77 | 63 |  |  |  |  |

Table 7 lists the statistical results of the chi-squared tests and correlation coefficients between the frequency values of the onset and codas of men and those of women in the Buckeye Corpus.

**Table 7.** Statistical results of chi-square tests and correlation coefficients between the frequency values of the onsets and codas of men and those of women in the Buckeye Corpus

|  | *df* | Chi-squared | *r* | *p* |
|---|---|---|---|---|
| Onsets | 92 | 1,200.5 | 0.983 | <.05* |
| Codas | 93 | 1,863.3 | 0.995 | <.05* |

The result of Pearson's chi-squared test between the onsets of the men and women was statistically significant ($\chi^2$=1,200.5, *df*=92, *p*<.05). On the other hand, the result of the chi-squared test between the codas of the men and women was also statistically significant ($\chi^2$=1,863.3, *df*=93, *p*<.05). We can conclude that there is a significant association between the phonological processes for both the onsets

and codas of men and women in the Buckeye Corpus. In addition, Pearson's product-moment correlation between the paired data was conducted to obtain $r=0.983$ ($p<.05$) for the onset data and nearly the same or stronger positive correlation coefficient for the codas, $r=0.995$ ($p<.05$). The two statistical results indicate that the phonological processes of the two sex groups correlate strongly with each other. It is quite interesting that the two sex groups show a proportionate ratio in their frequency distribution. Another noticeable trend is the deletion processes of men and women depending on the consonants. With /hh/, /r/, and /y/, men tend to delete more consonants than to preserve them compared to women. For example, 634 instances of /y/ of the onsets of men were deleted compared with 223 instances of /y/ of women. On the other hand, women tend to practice the same consonant processes more often than men. This result may reflect women's tendency to produce the sounds clearly as in the orthographic symbols. However, more cases of the same consonant processes for men than for women were also observed for the onsets /l, m, n, t, w/, a result that needs further study.

Here, we can extend the generalization on vowel processes (Yang, 2018) by including the current result on consonant processes: The distribution of the phonological processes of English vowels and consonants indicates a significantly strong association between men and women in the spontaneous speech corpus. As described in the previous study, this generalization may represent sex characteristics fairly well with sufficient data from the corpus.

Thus far, we have examined the general association between the data for men and women. If we go further and categorize those consonants by manner and place, we may find some specific association patterns, as in Kessler & Treiman (1997). They reported that the vowel-coda association was always stronger than the onset-vowel association after analyzing association patterns among the onsets, vowels, and codas of 2001 monomorphemic CVC words in the unabridged Random House Dictionary (Flexner, 1987). They added that the unequal associations indicate greater distinctiveness at the onset, which promotes more efficient and distinctive production and perception at the beginning of English words.

## 4. Summary and Conclusion

This study examined the phonological processes of English consonants in the spontaneous Buckeye Speech Corpus. R scripts were created to syllabify each word in both orthographic and pronounced forms. Then, the one-to-one phonological processes of English consonants were tabulated and divided into such syllable components as onsets and codas separately.

The results are as follows: First, the majority of the consonants were pronounced as the same sounds in conversation. The speakers in the Buckeye Corpus must have attempted to be more clearly understood. Second, the number of deletions in the codas is relatively greater than that in the onsets due to the information weight of the onsets. An observation of the total inventory distribution of consonants may have to be reexamined in light of the syllable components and their asymmetrical distributions of phonological processes. Third, we found that there is a significant association between the phonological processes of the onsets and codas in men and women in the Buckeye Corpus. Given these results, this paper concludes that an analysis of phonological processes in spontaneous speech corpora can improve the practical understanding of spoken English. Further studies would be desirable to compare the current

phonological process data with those of other languages to search for universal patterns in phonological processes.

This study could be applicable to future studies that extend not only to a specific language or dialect but also to a comparison of native and nonnative speech. For example, the phonological processes of native and nonnative speech may lead to interesting findings and applications, such as the establishment of better teaching plans or practices. Specifically, teachers may observe the general problems of students in a classroom and offer individually tailored practice lessons.

## References

Clements, G. (1990). The role of the sonority cycle in core syllabification. In J. Kingston, & M. Beckman (Eds.), *Papers in laboratory phonology 1: Between the grammar and physics of speech.* (pp. 283-333). Cambridge, UK: Cambridge University Press.

Cutler, A. (1982). The reliability of speech error data. In A. Cutler (Ed.), *Slips of the tongue and language production.* (pp. 7-28). Amsterdam, The Netherlands: De Gruyter Mouton.

Duanmu, S. (2002). Two theories of onset clusters. *Chinese Phonology*, *11*, 97-120.

Eckman, F. (2004). Universals, innateness and explanation in second language acquisition. *Studies in Language*, *28*, 682-703.

Ernestus, M., Hanique, I., & Verboom, E. (2015). The effect of speech situation on the occurrence of reduced word pronunciation variants. *Journal of Phonetics, 48*, 60-75.

Flexner, S. (Ed.) (1987). *The Random House dictionary of the English language.* New York, NY: Random House.

Kessler, B., & Treiman, R. (1997). Syllable structure and the distribution of phonemes in English syllables. *Journal of Memory and Language*, *37*, 295-311.

Maddieson, I. (1984). *Patterns of sounds.* Cambridge, UK: Cambridge University Press.

Noyer, R. (2016). Transcription of English syllable structure. Retrieved from http://www.ling.upenn.edu/~rnoyer/courses/103/Transcription.pdf/

Patterson, D., & Connine, C. (2001). Variant frequency in flap production: A corpus analysis of variant frequency in American English flap production. *Phonetica*, *58*, 254-275.

Pitt, M., Dilley, L., Johnson, K., Kiesling, S., Raymond, W., Hume, E., & Fosler-Lussier, E. (2007). Buckeye corpus of conversational speech (2nd release). Columbus, OH: Department of Psychology, Ohio State University. Retrieved from https://buckeyecorpus.osu.edu/

R Core Team. (2019). R: A language and environment for statistical computing (version 3.5.3) [Computer software]. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from https://www.R-project.org/

Williamson, G. (2014). Syllables and clusters. Retrieved from http://www.sltinfo.com/syllables-and-clusters/

Yang, B. (2012). Reduction and frequency analyses of vowels and consonants in the Buckeye Corpus. *Phonetics and Speech Sciences*, *4*(3), 75-83.

Yang, B. (2016). Phoneme distribution and syllable structure of entry words in the CMU English Pronouncing Dictionary. *Phonetics and Speech Sciences*, *8*(2), 11-16.

Yang, B. (2018). Phonological processes of vowels from orthographic to pronounced words in the Buckeye Corpus by sex and age

groups. *Phonetics and Speech Sciences*, *10*(2), 25-31.

Zue, V. & Laferriere, M. (1979). Acoustic study of medial /t, d/ in American English. *Journal of the Acoustical Society of America*, *66*(4), 1039-1050.

• **Byunggon Yang,** Corresponding author
English Education Dept.
Pusan National University
30 Changjundong, Keumjunggu
Pusan, 46261 Korea
Tel: 051-510-2619
Email: bgyang@pusan.ac.kr
Homepage: http://fonetiks.info/bgyang
Fields of interest: Phonetics, Phonology