

## *Syllable structure and word structure: a study of triconsonantal clusters in English*

JANET PIERREHUMBERT

### 11.1 Introduction

Current phonological theory maximizes the responsibility of the syllable for explaining co-occurrence restrictions on consonants. The inventory of word-initial consonant clusters is chiefly explained by the constraints on the syllable onset. The syllable coda also plays a central role in explaining which word-final consonants are permissible. The word node comes into play only by picking up extra peripheral elements, most notably the coronal appendices of English, and by defining the domain for any co-occurrence restrictions which cross syllable boundaries. (See Fudge 1969; Fujimura and Lovins 1978; Selkirk 1982; Clements and Keyser 1983.) With the phonotactic responsibility of the syllable thus maximized, the cross-product of codas and onsets is the starting point for any description of medial clusters. That is, in the absence of additional provisos, any concatenation of a well-formed coda and a well-formed onset is predicted to be possible medially in a word.

The present project evaluated the extent to which syllable structure explains the inventory of long medial clusters – the clusters of three or more consonants – in English. It was motivated by the observation that word-internally such clusters are extremely restricted in comparison with the set defined by the cross-product of codas and onsets. The basic model obviously requires modification, and a detailed examination of the occurring and missing clusters reveals what type of modification is needed.

Two methods were applied in the study. The pronunciation fields of the on-line Collins English Dictionary (distributed through the ACL Data Collection Initiative) were used to make an inventory of onsets, codas, and their frequencies, and to establish which medial clusters occur at all. Then, a

follow-up experimental study showed that members of the critical set of missing clusters represent systematic rather than accidental gaps.

The study deals only with clusters found medially in morphemes. That is, it deals with the triconsonantal clusters found intervocally in words such as “vanquish,” “lobster,” “doldrums,” “palfrey,” and “orchestra,” excluding those found in words such as “exactly,” “vastness,” “width-wise,” and “marksman.” Clusters found medially in bound morphemes were included, e.g. “anthro,” “andro.” However, the study does not deal with clusters occurring morpheme peripherally, such as /ntl/ in “gentl + er,” arising from /gentler/. Such cases present additional complications which we hope to investigate in a future study.

In the dictionary, 675 distinct clusters of three or more consonants are found. However, only fifty are found morpheme-medially (in a sense of “found” to be made more precise below). Compound words are by far the biggest source of long consonant clusters. The listing of compounds in the dictionary is of course spotty, and many more clusters would no doubt be found in a study of productively formed compounds.

In order to grasp the force of the number fifty, consider the number of different clusters which are taken to be well-formed according to the hypothesis that morphemes are arbitrary concatenations of syllables. The dictionary has 147 different consonantal sequences at the end of words and 129 at the beginning. (Words beginning or ending in a vowel are taken to represent a single case, that of no consonants.) Taking all possible combinations yields 18,963 possible medial sequences. This number is of course reduced by stripping appendices off final clusters and by enforcing a widely noted constraint against morpheme-initial geminates. Geminates are found only in compounds or across a word boundary, e.g. “subbasement.” Taking these generalizations into account only reduces the number of viable candidates to 8708. It’s a long way from 8708 to 50.

The assumption that the syllable grammar is stochastic was found to make the single greatest contribution towards addressing this discrepancy. The combination of a low-frequency coda and a low-frequency onset is expected to be a low-frequency occurrence. In fact, if the coda and following onset are statistically independent, then the probability of the combination is the product of the two low frequencies, and therefore far lower than that of either part. This means that many combinations are not expected to be found in a vocabulary of realistic size, even if both parts are found. It turned out that almost all occurring triconsonantal clusters were among the 200 most likely combinations, and that a stochastic interpretation of syllable grammar effectively ruled out a huge number of possible clusters, eliminating the need for many idiosyncratic constraints in the grammar.

However, it is still necessary to address the finding that only 50 of the 200 most likely combinations actually occur. Additional constraints, enforced at the morpheme or word level, are needed to rule out clusters which were likely a *priori*, but were not found. The experimental study established that these constraints represent part of the tacit knowledge of native speakers, rather than reflecting accidental gaps. The implications of the constraints for phonological theory are examined in section 11.6.

### 11.2 Methods I. Study of the dictionary

All arguments for hierarchical structure in linguistics have implicit *statistical assumptions*. We argue for constituents by showing that they serve as a domain for dependencies among elements. To defend a hierarchical structure, we must also show that elements which are equally close in the terminal string, but not claimed to be in the same constituent, do not exhibit such dependencies. If an equal degree of statistical dependence were found among all  $n$  adjacent elements, then no hierarchical structure would usefully distill the dependencies, and the data would suggest a quite different mathematical characterization, namely an  $n$ th-order Markov process.

The statistical viewpoint is particularly important in studies of the lexicon. The adult mental lexicon may be viewed as quasi-finite, with new forms added only slowly, and any given dictionary is certainly finite. Thus if a particular phonological combination is absent from the lexicon, it is necessary to establish whether its probability of occurrence is actually high enough that we would expect to find it in a sample the size of the lexicon or dictionary.

In order to address this issue, a method was adopted which was crude but nonetheless instructive. A pronouncing dictionary was extracted from the main dictionary by combining entries sharing both spelling and pronunciation, even if they differed in meaning. This was done because the dictionary uses polysemy to convey breadth of meaning; there are, for example, nine entries for "brother." Predicted probabilities for medial clusters of any length were then estimated by taking the cross-product, with frequencies, of all occurring word onsets and word-final syllable codas in the pronouncing dictionary. The medial clusters were then rank-ordered by predicted probability, from most to least likely. The remainder of the study then used the predicted probability rank (or the relative predicted probability), not the predicted probability itself.

The predicted probabilities were computed without regard to the morphological or etymological status of the words, for the sake both of

economy of effort (there are approximately 70,000 phonologically distinct entries in the dictionary) and of replicability. For example, no subjective judgment was made about which foreign borrowings are fully assimilated and which are not; it was assumed that the frequency of the consonant clusters in such words adequately represents their linguistic and cognitive status, with high-frequency combinations being fully acceptable even if they all came from the same donor language. Similarly, compounds were included on the grounds that all compounds have a good word beginning and a good word ending. The inclusion of compounds of course tends to inflate the contribution of words which are phonologically unusual but form part of many compounds. On this point, the crudeness of the approach may perhaps be excused by the results.

To tabulate coda frequencies, it is necessary to make a specific assumption about which word-final coronals are in the appendix and which are in the coda. For the present study, the conservative assumption is the one which maximizes the role of the appendix, thus minimizing the coda and accordingly minimizing the predicted number of medial clusters. The most conservative possible assumption would thus be the following:

- (1) Any word-final sequence of coronal obstruents is analyzed as being in the appendix.

According to this assumption, the appendix would cover not only the final /s/ in "pasts," but the entire /sts/ cluster. Indeed, it would cover the /t/ in "cat."

However, this extreme assumption cannot be maintained. It would make it impossible to syllabify the words in (2) (or whatever subset of these words the reader may judge to lack a word boundary within the cluster). The difficulty arises because word-final codas are being used as evidence about word-internal codas. Since (1) puts all coronals into an appendix rather than into a coda, it incorrectly implies that no word-internal coronal codas are permitted.

- (2)
- |           |              |
|-----------|--------------|
| vodka     | jodhpurs     |
| Atlanta   | pizza        |
| atlas     | Nazi         |
| badminton | bedlam       |
| chitling  | Presbyterian |
| jitney    | Aztec        |
| litmus    | husband      |
| nutmeg    | witness      |
| ordnance  | Frisbee      |
| apartment | antler       |

Therefore, a simple weakening of (1) was sought which would permit word-internal coronal codas. It was noted that in all the words in (2), the coronal follows a vowel, offglide, or nasal, never an obstruent or /l/. As a result, it was decided to count as codas only coronal obstruents which directly followed a vowel, offglide or nasal; those following any other phoneme were taken to be appendices. This means that the coronal in "cat" is put in the coda rather than in the appendix, and that in "vodka," etc., the coronal can also go in the coda. In "weft" the /t/ is in the appendix and the /f/ counted towards the tally for /f/ in coda position. This treatment of the appendix still yields a rather conservative estimate of the number of possible medial clusters.<sup>1</sup>

The constraint against geminates was not enforced in constructing the list of possible clusters. There was no reason to do so since the treatment of geminates does not affect the way in which nongeminate clusters rank with respect to each other in expected probability. Including clusters with geminates in the list provides an opportunity to compare the statistical behaviour of a known constraint to the behavior of constraints emerging from the study.

Once the rank-ordered list of properly syllabifiable clusters of three or more elements was constructed, it was then compared to a list of clusters of three or more consonants which actually occur morpheme-medially. The latter list was constructed by extracting all words in the dictionary with three or more consonants in a row in the pronunciation field and sorting by the cluster exemplified. The set of entries for each cluster was then read to determine if it included any in which the cluster was morpheme-medial.

A few comments are in order about this determination, clearly the most subjective step in the entire process. A cluster was taken to occur if it occurred morpheme-medially in at least two reasonably familiar words. No effort was made to establish or interpret rates of occurrence, which indeed varied widely and not always as predicted. The class of "reasonably familiar" words was taken to include words such as: "pancreas," "extirpate," "palfry," "doldrums," "velcro," "imbroglio," "eclampsia," "wainscot." Examples of words in the dictionary which were not taken to be "reasonably familiar" include: "Melanchthon," "rigsdaler," "hoactzin," "anschluss," "pozzuolana." "Monomorphemic words" were taken to include a number of Greco-Latinate words which are historically polymorphemic, but which, it was felt, were probably not decomposed by most present day speakers. Examples include: "complete," "extreme," "inspect," "obtuse." Words such as "exhusband" and "Transsiberian" were of course taken to be polymorphemic, as were words in which the meaning

of an affix was discernible even if the meaning of the entire word could not be determined compositionally. Examples of the latter type include "excavate," "exclude," "excrete," and "excursion" (all sharing a meaning of "ex" as "out") and "transparent," "transform," and "transfer" (sharing the use of "trans" to indicate change of location or state).

Many generative linguists would decompose "complete," "extreme," "inspect," and "obtuse," following principles laid out in Nida (1949). These principles permit the isolation of meaningless morphemes, or formatives, provided that they form the residue when an independently meaningful part is taken away. For example, "con" in "condense" can be isolated because "dense" is independently found with the appropriate meaning. Then, a principle of transitive closure on isolatability permits the isolation of "con" in "condense" to support the isolation of e.g. "flict" in "conflict," "trol" in "control," and so on. When these principles are applied in a literal fashion, they lead to ludicrous overdecomposition. For example, by isolating "con," we get "con + quer," leading to "bi + cker," "han + ker" and "pu + cker"; similarly the noun "con + tra" supports the decomposition "Char + tres" (according to the British pronunciation listed in the dictionary) and the decompositions "king + dom" and "con + dom" support "a + dam," "ma + dam" and "maca + dam." The fact that such decompositions have not been proposed in practice suggests that scholars have implicitly applied their knowledge of semantics, spelling, and historical development. There's no reason to suppose that the intuitions of people with so much linguistic sophistication and training would be shared by the ideal naive speaker-listener. Psycholinguists have in general been far more conservative about assuming that forms involving semantically opaque and nonproductive derivational morphology are synchronically decomposed; see, for example Bradley (1979), Bybee (1988), Nagy and Anderson (1984). As a result, I am inclined to agree with the position expressed in Bybee (1988), according to which the identification of a meaningful subpart of a word does not imply that the residue is also a morpheme.

A number of possibly interesting points were not addressed in the study. Since the dictionary has British pronunciations, there are no postvocalic /r/s and any questions concerning their behavior in American English cannot be addressed. The palatal onglide (as in "tune") was not treated as a consonant. The possible role of stress in conditioning medial clusters was not investigated. Due to the large number of noun-verb pairs differing only in stress (e.g. "conflict," "con'flict") it was judged that stress would not be a primary influence on the form of medial clusters. However, the possible role of stress deserves further attention, in particular the relation of stress to homorganicity requirements for nasals.

### 11.3 Results of the dictionary study

The assumption that the syllable grammar is stochastic, with the likelihood of medial clusters derived from the independent likelihoods of the component codas and onsets, made an extremely successful contribution to the characterization of medial clusters. This success is displayed in figure 11.1. To construct this figure, the triconsonantal clusters as ordered by likelihood were arbitrarily grouped in twenties: "1" on the x-axis represents the group of the most likely twenty, "2" represents the group of the next most likely twenty (that is, clusters ranked 21-40), and so on. The y-axis shows how many in each group are actually found.

The top ten groups (or the 200 most likely clusters) include practically all those found; a single group of exceptions will be discussed below. The predicted rate of occurrence for the 200th cluster is approximately 1 in 10,000. Though the method used did not actually provide a count of the number of polysyllabic monomorphemic words, it may be noted that the dictionary had about 70,000 distinct entries, with a very large number of these words being polymorphemic or monosyllabic. Thus, the cutoff has a realistic relationship to the size of the dictionary. The figure also shows that the rate of occurrence decreases as the predicted likelihood goes down.

Figure 11.1 also shows that a stochastic syllable grammar is not the whole story. It reduces the number of candidates to 200, but (as already noted)

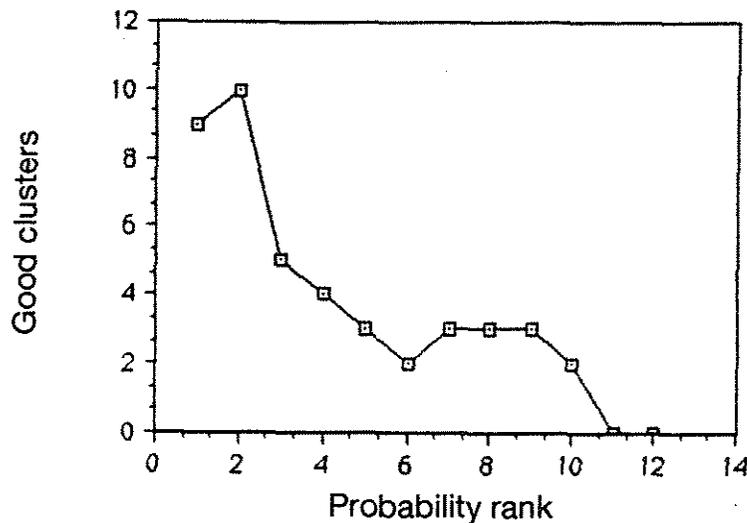


Figure 11.1. Occurring clusters per twenty candidates.

only about fifty are found. Even for the forty most likely candidates (all with expected rates of occurrence above 1 in 1842), just fewer than half are found.

Examination of the occurring and nonoccurring clusters in the top 200 candidates revealed the following generalizations:

First, nasal-stop sequences agree in labiality: either both phonemes are labial, or neither. This generalization covers eleven clusters which would otherwise be sufficiently probable that one would expect to find examples of them in the dictionary. Of these eleven, four involve /n/ before a labial; five have /m/ before a velar; and two have /m/ before a coronal. None of these in fact occurs. However, agreement in labiality is not enforced before fricatives, or at least not as strictly. /mst/ and /msp/ are marginally attested: /n/ does occur before /f/.

Second, clusters with a coronal obstruent in the coda do not occur.<sup>2</sup>A total of 79 cases are covered by this generalization (excluding clusters with geminates). This is by far the strongest generalization observed. Thirty-seven cases covered involve /t/ or /d/ in the second coda position following a nasal. 42 involve /t/, /d/, or /s/ preceding a biconsonantal onset. Triconsonantal clusters such as /str/ which comprise a well-formed onset are of course allowed.

Reference to the biconsonantal clusters found in the study of the appendix suggests a way to understand this generalization. As noted above, coronal obstruents must be permissible in coda position in order to syllabify occurring biconsonantal sequences. However, cases of biconsonantal medial clusters with a coronal obstruent in coda position are obviously rare compared to what one would expect from the overall phoneme frequencies for /t/, /s/, and /d/. This pattern is partly explained by the fact that final coronals in Latin prefixes have assimilated in place of articulation to following obstruents. (See e.g. Nesfield 1898: §570). In fact such assimilation of coronals is very common and probably has occurred in some of the many other languages which have contributed to the English vocabulary. Speaking synchronically rather than diachronically, the observation is that the coronal obstruents have a lowered probability internally as opposed to word-finally. If the probabilities for the triconsonantal clusters were estimated using the internal rather than the word-final frequencies, then triconsonantal clusters with coronal coda obstruents would fall below the threshold for occurrence. The linguistic problem thus reduces to that of enforcing position-dependent probabilities.

Third, velar obstruents occurred only before coronals in the clusters studied, never before labials or other velars. This generalization covers 14 clusters which would otherwise be expected to occur. Seven of the clusters ruled out involve a /k/ before a labial. Five involve /k/ before another velar.

Of these, three are independently covered by antigemination, and the remaining two would be covered if antigemination were to disregard voicing. Two involve the velar nasal before a labial. The velar nasal (having been tabulated word-finally) is interpreted phonologically as /ng/ before a labial, where it could not have arisen from assimilation. It is interesting to note that /ng/ is actually more common as a coda than /g/ alone. No clusters beginning with /g/ were sufficiently likely to be in the running.

It may be noted that velar obstruents do occur before labials in biconsonantal clusters: "rigmarole," "pigmy," "dogma," "Egbert," "rugby," "tacmahack," "acme," "Micmac," "Achmed." Furthermore, in the word "angma" (not in the dictionary), an underlying /g/ is presumed to precede a labial in a triconsonantal cluster, although this /g/ does not appear as an obstruent on the surface. However, velars before noncoronals appear unexpectedly rare. This fact could be described in terms of position-dependent effects on probability, just as the shortage of coda coronal obstruents was.

Fourth, as expected, there were no clusters involving geminates. This generalization covers seventeen cases; however, of these, twelve also had unacceptable coronal obstruents in coda position.

Fifth, in addition to the lack of geminates, a lack of clusters with identical first and third elements was also observed. Clusters of this form falling in the most likely 200 are:

- (3) /fl/ /kl/ /pl/ /bl/ /gl/ /sl/  
 /tst/ /ntn/ /ndn/ /tstr/ /nsn/  
 /ksk/

Of these, four are also excluded by the restriction against coronal-obstruent codas. /nsn/ can be syllabified with /s/ in the onset. All but two involve /l/, already observed by Clements and Keyser (1983) to be subject to a dissimilarity requirement between the onset and the coda. Therefore the status of /nsn/ and /ksk/ is of particular interest. /ksk/ actually occurs in a number of words beginning with "ex." Of these, quite a number are viewed here as decomposable because they contain the meaning element usually associated with "ex"; in many cases the decomposition is further supported by related forms. Such cases include "excommunicate" (cf. communicate), "exculpate" (cf. culpable), "exclaim" (cf. claim, disclaim, declaim), "excavate" (cf. cavity), "excruciate" (cf. crucifix). However in three cases ("excuse," "Excaliber," and "exquisite"), the support for decomposition is less apparent. The unclear status of these examples led to /ksk/ being included in the experimental study, so the issue will be taken up again below.

Enforcing these five generalizations leaves us with the following unexplained gaps:

- (4) /lpr/ /nsm/ /lbr/ /nsw/ /lsp/ /lkw/ /kdr/ /ngtr/ /ndw/ /ksl/  
 /vpr/ /ksw/ /ngks/ /ptr/ /nspr/ /lstr/ /pst/ /vtr/ /vst/ /pkr/ /mstr/

Of these clusters, the following occur marginally in single examples, partially unassimilated borrowings, or proper names which may be decomposed:

- (5) /lpr/ culprit  
 /lbr/ Galbraith, Albrecht  
 /nsw/ mansuetude, consuetude  
 /lsp/ felspar  
 /lkw/ Alcuin  
 /ngtr/ Langtry  
 /ndw/ Gondwana  
 /ksl/ Bexley, Huxley  
 /ksw/ Maxwell  
 /ngks/ Yangtze; common word-finally, e.g. Bronx  
 /ptr/ calyptra  
 /lstr/ maelstrom  
 /pst/ capstan, Epstein  
 /mstr/ Armstrong

The reader is left to his or her own conclusions about the status of these examples. /nsm/ and /nspr/ occur only in compounds and across word boundaries (advancement, mainspring, etc.).

The following are completely absent from the dictionary.

- (6) /kdr/, /vpr/, /vtr/, /vst/, /pkr/

Since no clusters beginning in /v/ are good, one might propose a frequency dependence on position for /v/ as for /t/. However, with only three relevant examples, all of which are already none too likely, it is difficult to say whether this gap is accidental or systematic.

Study of the dictionary also revealed a group of examples which are found in defiance of their unlikelihood. These are clusters beginning in /b/ followed by an /s/ or /t/, almost all historically originating from the prefixes "sub-," "ob-," and "ab-." Examples are:

- (7) /btr/ subtract, obtrude  
 /bst/ abstain, substance, lobster, obstetric, obstinate, substitute  
 /bsk/ obscure  
 /bskr/ subscribe  
 /bstr/ abstruse, abstract, obstruent, obstruct, obstreperous

The expected frequencies for these combinations are in the range of 1/20,000 to 1/80,000.

Any mechanism which can decrease the likelihood of /t/ morpheme-internally can also be applied to increasing the likelihood of /b/. In this sense, the examples are not problematic. However, they raise the issue of the extent to which these clusters are acceptable in other phonemic contexts than they usually appear in. Is "tibstance" a possible word of English? How about "chabtry"? If the acceptability of the clusters depends on their broader phonological form, this would tend to support the idea that the degree of overall similarity between a word and the others in the lexicon determines how good it is phonotactically. Suggestions of this sort have been made by Greenberg and Jenkins (1964) and Ohala and Ohala (1986).

#### 11.4 Methods II: experiment

In order to verify that a number of the observed constraints actually represent aspects of the tacit knowledge of native speakers, a small experiment was carried out.

Sixteen actually occurring clusters were selected which had a range of predicted frequencies and which are found in at least two reasonably familiar disyllabic words. These clusters were:

- (8) /str/ /spr/ /skr/  
 /ntr/ /nst/ /nkr/ /ngr/ /nkw/ /nkl/ /ndr/  
 /mpr/ /mbr/ /mpl/  
 /lst/ /lkr/  
 /kst/

Note that the set of occurring clusters is highly unbalanced phonologically, and thus it was impossible to select a phonologically balanced experimental set. An effort was made for a reasonable degree of diversity. However, the possibility of artifactual effects due to phoneme imbalance within the experiment (as opposed to within the English language) cannot be excluded.

Sixteen bad clusters were also selected. They all had predicted frequencies within the range for good clusters, and violated one of the above observations. The bad clusters used were:

- (9) /fl/ /kl/ /ksk/ /nsn/ (duplicate consonants)  
 /tkl/ /dbr/ /dgr/ /tpr/ /tfl/ (coronal coda)  
 /ʃp/ /ʃk/ /ʃm/ (coronal coda)  
 /mkr/ /mgr/ /mtr/ (/m/ before a nonlabial stop)  
 /mst/ (/m/ before a nonlabial obstruent)

Clusters with /ʃ/ were included to provide examples in which a coronal obstruent was clearly not a morphologically separate appendix. They were the only biconsonantal clusters in the experiment; triconsonantal clusters with /ʃ/ had predicted frequencies too low to be included. /mst/ was included as an example of a cluster which is only weakly unacceptable. It provides a baseline for the results for the other clusters.

Forty-eight existing disyllabic words (actually containing one of the good clusters) were used to create 48 "good" or simplex nonsense words and 48 "bad" or compound nonsense words. The nonsense words were created by substituting at random a different medial cluster for the one that actually occurred. After this random substitution, further switches were made to remove substrings corresponding to actual words as actually spelled. Particularly recalcitrant examples which made it impossible to remove actual words while maintaining the balance of the set were resolved by altering a consonant in the base word to create a new base form. This new base was used for both the "good" and the "bad" versions.

This procedure resulted in a "good" word and a "bad" word formed from each base, e.g.

- |      |          |          |          |
|------|----------|----------|----------|
| (10) | BASE     | GOOD     | BAD      |
|      | bistro   | bimplo   | bilflo   |
|      | constant | cosprant | comkrant |

Each good cluster was found in three different words, paired with a variety of bad clusters. Similarly, each bad cluster was found in three words, paired with a variety of good clusters.

Words were presented in ordinary English spelling to an undergraduate linguistics class. Nasals preceding a velar were written "n," even if presumed to be homorganic, e.g.:

- (11) tancrum, pongrete

/kw/ was spelled with a "qu":

- (12) fonquess, inquisite

/ks/ was spelled with an "x":

- (13) traxtil, uxkage

The palato-alveolar fricative was written "sh," somewhat disguising its anomalous status in the experiment.

Ordinary spelling was used because many of the undergraduates had a poor grasp of transcription. The course in question was an introductory class satisfying a distribution requirement and covering only basic concepts of phonology and phonetics. As a result, the subjects can be assumed to

have had at most a very general understanding of the aims of the experiment.

The students were told that the words were candidates for the vocabulary of a science-fiction novel. Half were asked to judge which word of a matched pair was "most like a compound." The other half were asked to judge which seemed "most suitable to be part of the vocabulary of an English speaker of the 21st century." All subjects were instructed to work quickly, giving their first impression, and to answer all questions even if they had to guess. There were twelve sets of responses for each set of instructions. Subjects working under one set of instructions were not aware that others had a different set of instructions.

The presentation randomized the position (first or second) of the "good" cluster of each pair. The different pairs were also randomized with each other. Matched pairs (e.g. "cosprant," "comkrant") were presented rather than unmatched pairs (e.g. "bimplo," "comkrant") in order to facilitate the comparison. Calling attention in this way to the contrast under investigation obviously maximizes the chance of finding a difference. It was felt that success with this format of presentation would lay the groundwork for a more elaborate experiment disguising the contrast under study. That is, if the present design failed to produce results there would be no point in continuing. However, the present experiment must be viewed as a pilot and a full-scale study should also be carried out.

### 11.5 Results of the experiment

The responses to both sets of instructions show that subjects had tacit knowledge of the regularities in triconsonantal clusters and could apply this knowledge in evaluating novel forms.

For the instruction "which is a compound," 40 out of 48 possible word comparisons went in the direction predicted (that is, more than six subjects selected the intended compound as a compound.) There were three ties and five comparisons were contrary to prediction. Scores were tabulated for each cluster by combining scores for all three pairs in which it occurred. Of the 16 clusters predicted to occur only in compounds, 15 were judged to occur in compounds more than half the time. There was a tie for /mst/, the cluster which had been included as only marginally problematic.

For the instruction "which is more suitable," 42 word comparisons came out as predicted with three ties and two contrary to prediction. All 16 bad clusters were judged to be "more suitable" less than half the time.

Tabulating by groups yields the following results. Numbers represent percents of judgments pooled across subjects and words.

(14)

|                              | Compound | Suitable |
|------------------------------|----------|----------|
| <i>Coronal in coda:</i>      | 71       | 22       |
| /ʃ/ in coda:                 | 76       | 19       |
| /t/ or /d/:                  | 68       | 23       |
| <i>Duplicate consonants:</i> | 66       | 20       |
| /ksk/ alone                  | 75       | 17       |
| <i>Nonhomorganic /m/:</i>    | 58       | 26       |
| excluding /mst/:             | 61       | 19       |

It is interesting to note that the constraints against coronals in the coda and against duplicate consonants are both stronger than the tendency to avoid non-homorganic nasals, which has been previously noted in the literature. This result was obtained even using the nonhomorganic clusters of clearest status, those involving /m/ before a nonhomorganic stop. It may also be noted that /ksk/ was the least acceptable of the clusters involving duplicate consonants, and the most likely to be viewed as arising from a compound. This result tends to support the claim that the constraint is not idiosyncratic to /l/. However, further work is needed to rule out the possibility that this result is an artifactual result of spelling with an "x."

In interpreting the results, it is necessary to rule out the possibility that they are adequately explained by differences in predicted frequency between the good and bad clusters. This possibility must be evaluated, because the experimental design did not actually control for predicted frequency. No bad cluster had a predicted frequency as high as the most likely good clusters. The materials were designed in this way because the aims in designing the materials were to some extent at odds with each other. These aims were: to use several different clusters to evaluate each proposed constraint; to control for predicted frequency; and to include some clusters which were unequivocally acceptable, in view of the actual rarity of some occurring clusters. Because of imbalance in predicted frequencies for good and bad clusters, pooled data would be expected to exhibit some tendencies in the direction noted even under the null hypothesis. If subjects simply selected the most probable cluster of each pair, the judgments would on the average favor the "good" clusters.

A subset of the "compound" data was extracted in order to eliminate this possibility. The ten pairs in which the predicted frequency of the "bad" cluster strictly exceeded the predicted frequency of the "good" cluster were extracted. Out of 120 individual judgments, 81 were nonetheless as predicted. This is highly significant by a binomial test. Pooled data for

eight word pairs were as predicted, with one tie and one contrary to prediction. It may be concluded that subjects do have phonotactic knowledge of the triconsonantal clusters, apart from that implied by a stochastic syllable grammar.

### 11.6 Discussion and conclusions

A stochastic model of syllable structure goes far towards explaining which triconsonantal clusters are found. The extent to which the clusters can be generated as statistically independent selections of a coda and a following onset confirms the existence of the syllable as a unit of hierarchical structure. It provides evidence against the view that the form of medial clusters is determined entirely by sequential constraints. The form of the evidence is brought out by considering the alternative, that sequential constraints in the form of a finite state model define the allowable medial clusters. Under this approach, a separate model is needed for medial clusters since many occur neither initially nor finally. The fact that the statistics of initial and final clusters so effectively circumscribe the medial alternatives is treated as accidental under this approach.

Statistical knowledge of phonological structure has also been demonstrated in two other areas. Experiments described in Kelly (1988) show that English speakers are able to apply statistical knowledge of the rhythmic contexts for nouns and verbs. Cassidy and Kelly (1991) show that English speakers are aware that nouns are typically longer than verbs.

One of the observed constraints above the syllable level can be adequately described in current phonological theory using marking conditions. Condition (15) prevents a nasal from preceding a stop which does not agree in the feature [labial]:

- (15)    \* [nasal]    [-cont]  
           |            |  
           C            C  
           |            |  
       [α labial]    [-α labial]

If [labial] is viewed as a privative feature, then it becomes difficult to collapse the cases of e.g. \*/mk/ and \*/np/.

A reviewer suggests that the generalization stated in (15) be attributed to failure of the coda to license a place of articulation for nasals. Place of articulation would either be acquired by a spreading rule from the following stop, or default to coronal. The marginal existence of a contrast between /m/ and /n/ before fricatives provides equally marginal counterevidence to this

suggestion. In addition, under this approach it is unclear how to cope with the fair number of English words containing coda /m/ before onset /n/.

Let us now turn to the cases which do have implications for autosegmental theory. Consider first the distribution of coronal obstruents in coda position. Such obstruents are found to be far more frequent in word-final position than medially; that is, word and syllable position interact to determine frequency. This frequency difference (which is supported by a study of biconsonantal clusters) effectively predicts the absence of triconsonantal clusters beginning with coronal obstruents. Similarly, velars are found to be unexpectedly infrequent in coda position preceding noncoronals. Their infrequency in this position, as revealed by their sporadic occurrence in biconsonantal clusters, effectively precludes them from appearing in the triconsonantal clusters.

The central idea of current licensing theory is that structural nodes (such as the syllable, the onset, or the coda) support phonological contrasts by supporting segmental features. If a particular structural position – for example the coda – does not display the full set of contrasts which are available in the language, this is because that particular structural position does not license all the features in the language. Licensing theory provides two examples of how nodes can interact to control segmental content. The first is the treatment of the appendix. The syllable coda can license a single stop, which can have any place of articulation. The word node licenses additional segments following the coda (that is, the appendix), which must be coronal. These segments might be direct structural dependents of the word node. Or else they could be structural dependents of the coda which are only permissible when this coda is word final, as argued in Scobbie (1991). But in either case, the syllable licenses all the different place features but the word node doesn't. The interaction between the syllable and the word is seen from the fact that the segments licensed by the word node come after those licensed by the syllable.

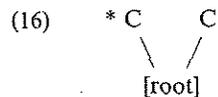
It is also possible for two nodes to jointly control different aspects of the same segment. This is the case of parasitic licensing, as studied in Itô (1986). In a number of languages, including Japanese, a consonant is only permitted in the coda if it agrees in place with the onset of the next syllable. That is, the only distinctive feature in this position is [nasal], with the outcome being either a nasal homorganic cluster or a geminate. This situation is described by permitting the coda to license only [nasal], with the place features formally originating with the following syllable onset. Itô's treatment relies on the use of negative marking conditions with links interpreted exhaustively; however, the same idea can be carried out using positive conditions, as described in Goldsmith (1990) and Scobbie (1991).

In cases of parasitic licensing studied to date, the two nodes controlling features of a segment are adjacent syllable nodes. However, nothing in principle prevents hierarchically related nodes from behaving in the same way. Note however, that the pattern found here for coronals and velars in coda position involves two nodes jointly controlling the same set of features, not merely different features of the same phoneme. The situation thus requires a straightforward formal extension of the present formalism, permitting hierarchically superior nodes to readjust or load probabilities assigned lower down in the hierarchy. In the case of the velars, this readjustment must also refer to the onset of the following syllable whose place of articulation is relevant. In general, this type of interaction among nodes is only brought out in a statistical treatment of phonotactics.

It is interesting that a statistical effect operating above the syllable level is found to effectively describe two absolute patterns in the triconsonantal clusters. It would of course be possible to formulate a purely qualitative marking condition which prohibited e.g. coronal codas in triconsonantal clusters only. However, this description would fail to relate the situation in triconsonantal clusters to that in biconsonantal clusters and would thus be less succinct and general than the description proposed here. These two cases thus raise an important issue, namely to what extent the notion of an "ill-formed" word can be reduced to that of a "statistically improbable" word.

The required dissimilarity between the first and third consonants in a cluster is also interesting. There are two relevant precedents in the literature, but neither of them can account for this dissimilarity.

To prevent the occurrence of morpheme internal geminates, English is taken to have a marking condition of the following form:



(That is, two adjacent consonants cannot share all features). However, this condition does not preclude identity of the first and third consonants of a cluster, because they are not adjacent; the second consonant intervenes.

There is also by now a substantial literature on dissimilarity requirements which affect nonadjacent as well as adjacent consonants. Effects on nonadjacent consonants are described using the combined assumptions of autosegmental projection, the Obligatory Contour Principle (OCP), and underspecification theory. The OCP says that two adjacent like elements are prohibited. This idea, introduced in work on tone by Leben (1973) and christened in Goldsmith (1976), was then applied in McCarthy (1986) to problems in segmental morphophonology in Arabic and other languages.

By applying the OCP only to some particular featural projection, it is possible to rule out sequences of segments which are similar in some particular respect without being completely identical. For example, disallowing two identical specifications on the nasal tier would make it impossible to have two nasals of any type in a row. (It would also rule out sequences which agreed in being nonnasal, without the further understanding that nonnasal segments are unspecified for nasality rather than being [-nasal].) Analyses of this sort are Yip's (1988) treatment of voicing restrictions in Japanese, Steriade's (1987) treatment of /r/ - /l/ alternation in Latin, and McCarthy's description (this volume) of cooccurrence restrictions in the verbal roots of Arabic. These analyses all exploit underspecification and/or privativity to make certain segments transparent to OCP effects. That is, if a segment is not specified for some feature, then it will be invisible on the tier for that feature, and otherwise nonadjacent segments will be rendered effectively adjacent. In combination, then, the OCP and underspecification or privativity can effectively prevent even nonadjacent phonemes from sharing some phonological properties.

In attempting to apply this approach to the present problem, the absence of clusters /ksk/, /nsn/, /lll/, /lpl/, and /lbl/ turns out to be particularly important. The first two undercut the otherwise plausible suggestion that the gaps are explained by an OCP effect on the [lateral] tier, since the first and third consonants are not lateral. The last three, involving a medial labial, show that underspecification for coronals, as advocated in Paradis and Prunet (1991), cannot be exploited to render the first and third consonants effectively adjacent (either on the place tier or with respect to one of the place features). Furthermore, if coronals were unspecified for place, then not only /s/ but also /n/ would be unspecified. As a result, /nsn/ could not represent an OCP violation with respect to place, but only with respect to [nasal] or [continuant]. However, these features cannot be the domain for the constraint because of the contrast between acceptable and nonacceptable clusters which are entirely nonnasal, and because of the contrast between \*/nsn/ and /nst/, /nsk/, /nsp/.

In general, it is impossible for the OCP to prevent total identity across arbitrary intervening material, if its operation is restricted to features which are strictly adjacent on a tier. This is impossible because checking for total identity requires examination of all the tiers. But some of these tiers will be tiers on which the features of intervening material appear. If two segments are not adjacent, then there is some tier on which their features are not adjacent, for some choice of intervening material.

In view of this difficulty, the present study was followed up with a statistical study of the verbal roots of Arabic (Pierrehumbert 1993). These roots consist of three consonants, with dissimilarity requirements affecting

nonadjacent as well as adjacent consonants (see McCarthy, this volume). They thus provide a very relevant comparison to the medial clusters of English; but because they are so much more numerous and varied, they provide more detailed evidence about the formal character of dissimilarity requirements. The study, completed just as the present article goes to press, demonstrates that even in Arabic the dissimilarity constraints must be permitted to refer to nonadjacent feature specifications. The constraint against total identity is shown to be more persistent (better able to cross intervening material) than the constraint against mere homorganicity, and a model derived from the psychological literature which exhibits this behavior is laid out.

The lack of clusters in English with identical first and third consonants is thus consistent with a more extensive pattern in Arabic, the now classic example of a language with OCP effects on segments. We conclude, therefore, that English provides a further example of a dissimilarity requirement operating across intervening material. This conclusion is already anticipated by Clements and Keyser (1983), who note the absence of words like "fill," and Davis (1989), who established the systematic absence of sequences such as "ssep." Future work will need to establish why words in which one of two identical consonants is word-initial (such as "lilt" and "cake") are apparently exempted from such a constraint.

#### Notes

- 1 With hindsight, the author would suggest that postnasal coronals (as in /pænt/) probably actually count as appendices. The main consequence of this treatment would have been to raise the already high frequency tally for coda /n/, and to absolutely preclude internal clusters in which /n/ precedes a coronal obstruent that cannot be syllabified in the following onset. In short, "antler" and "handsome" would necessarily be polymorphemic under this view. We also observe that the possibility of a long nucleus before a coronal obstruent coda hangs by the thread of the word "ordnance" (pronounced with a long vowel in British English and a rhotic offglide in many American dialects).
- 2 In making this generalization /ntl/ is taken to be nonoccurring even though it is actually found in "antler." This is because the word is the sole exemplar of this cluster, once having set aside words in which the entire cluster is morpheme-final. The acceptability of "antler" may be related to the large number of such words in which a morpheme-final /l/ is syllabified with the following suffix.

#### References

- Bradley, D.C. 1979. Lexical representation of derivational relations. In M. Aronoff and M.L. Kean (eds.) *Juncture*. Cambridge, MA: MIT Press, 37-55.

- Bybee, J. 1988. Morphology as lexical organization. In M. Hammond and M. Noon (eds.) *Theoretical Morphology*. New York: Academic Press, 119-142.
- Cassidy, K.W. and M.H. Kelly. 1991. Phonological information for grammatical category assignments. *Journal of Memory and Language* 30: 348-369.
- Clements, G.N. and S.J. Keyser. 1983. *CV Phonology: A Generative Theory of the Syllable*. Cambridge, MA: MIT Press.
- Davis, Stuart. 1989. Cross-vowel phonotactic constraints. *Computational Linguistics* 15(2): 109-111.
- Fudge, E.C. 1969. Syllables. *Journal of Linguistics* 5: 253-286.
- Fujimura, O. 1990. Demisyllables as sets of features. In J. Kingston and M. Beckman (eds.) *Papers in Laboratory Phonology I. Between the Grammar and Physics of Speech*. Cambridge: University Press.
- Fujimura, O. and J. Lovins. 1978. Syllables as concatenative phonetic units. In Bell and Hooper (eds.) *Syllables and Segments*. Amsterdam: North-Holland.
- Goldsmith, J. 1976. Autosegmental phonology. Ph.D. dissertation, MIT. Published by Garland Press, New York, 1979.
1990. *Autosegmental and Metrical Phonology*. Oxford: Blackwell.
- Greenberg, J. H. and J. J. Jenkins. 1964. Studies in the psychological correlates of the sound system of American English. *Word* 20: 157-177.
- Itô, J. 1986. Syllable theory in prosodic phonology. Ph.D. dissertation, University of Massachusetts at Amherst. Published by Garland Press, New York, 1988.
- Kelly, M. H. 1988. Rhythmic alternation and lexical stress differences in English. *Cognition* 30: 107-137.
- Leben, W. R. 1973. Suprasegmental phonology. Ph.D. dissertation, MIT. Published by Garland Press, New York.
- McCarthy, J. 1986. OCP Effects: gemination and antigemination. *Linguistic Inquiry* 17: 207-265.
- Nagy, W. and R. Anderson. 1984. How many words are there in printed school English? *Reading Research Quarterly* 19: 304-330.
- Nesfield, J. C. 1898. *English Grammar Past and Present*. New York and London: Macmillan and Co.
- Nida, E. 1949. *Morphology: The Descriptive Analysis of Words*. Ann Arbor: The University of Michigan Press.
- Ohala, J. J. and M. Ohala. 1986. Testing hypotheses regarding the psychological manifestation of morpheme structure constraints. In J. J. Ohala and J. J. Jaeger (eds.) *Experimental Phonology*. Orlando: Academic Press, 239-252.
- Paradis, C. and J-F. Prunet. 1991. *The Special Status of Coronals: Internal and External Evidence*. San Diego: Academic Press.
- Pierrehumbert, J. 1993. Dissimilarity in the Arabic verbal roots. In NELS 23. Amherst, MA: University of Massachusetts at Amherst, GLSA Publications.
- Scobbie, J.M. 1991. Attribute value phonology. Ph.D. dissertation, University of Edinburgh.
- Selkirk, E.O. 1982. The syllable. In H. van der Hulst and N. Smith (eds.) *The Structure of Phonological Representations*, Parts I and II. Dordrecht: Foris.

*Syllables*

- Steriade, D. 1987. Redundant values. In A. Bosch, B. Need, and E. Schiller (eds.) *Papers from the Parasession on Metrical and Autosegmental Phonology (23rd Regional Meeting of the Chicago Linguistic Society)*, 339–362.
- Yip, Moira. 1988. The obligatory contour principle: a loss of identity. *Linguistic Inquiry* 19: 65–100.

***Part III***  
***Feature Theory***

---