

Structural Definition of Affixes from Multisyllable Words

by Lois L. Earl,* Lockheed Missiles and Space Company, Palo Alto, California

In a recent paper by H. L. Resnikoff and J. L. Dolby, "The Nature of Affixing in Written English," an algorithm for the structural definition of affixes was developed and applied to data consisting of all the words of the form CVCVC in the Shorter Oxford Dictionary. Fourteen strong prefixes and twelve strong suffixes and seven weak prefixes and forty weak suffixes were defined, but it was noted that all the affixes could not be expected to show up in two-vowel-string words. This paper summarizes the results of applying a modified form of the operational definition to data consisting of all the four-, five-, six-, and seven-vowel-string words in Webster's Third New International Dictionary. Thirteen additional weak suffixes, nineteen weak prefixes, seventeen strong prefixes, one strong suffix, and twelve possible suffix-compounding elements were found.

In this paper, as in the preceding one,¹ the aim is to define affixes from structural criteria alone. The problem of when an affix sequence is genuinely acting as an affix (as *re* may be considered a prefix in *react* but not in *read*) will not be considered, though the categorization into strong and weak affixes is intended to anticipate this problem. The validity of the defined affixes will be indicated only by comparison with existent affix lists. A more utilitarian evaluation of their validity can be made after the syntactic and phonetic implications of the defined affixes have been investigated.

The definitions for affixes given in this paper are essentially unchanged but are extended to include both one- and two-syllable affixes. The data set to which these definitions are applied is the four-, five-, six-, and seven-vowel-string words, a set of about 11,250 words. From this set the one-vowel-string affixes that did not occur in the two-vowel-string data set (used in reference one) will be defined, along with the two-vowel-string affixes that could not have occurred in the two-vowel-string data.

The extended definition for strong prefixes can be summarized as follows (consonant strings referred to in the definition are given in Table 1): Given a word of the form $C_1V_1C_2V_2C_3V_3 \dots$, if either C_2 or C_3 is an inadmissible consonant string, there is a mandatory syllabic break within the string, and everything preceding that break is defined as a "prefix possibility." A prefix possibility is defined as a "prefix probability" if in the data there are at least four words with the same prefix possibility arising from the same consonant string. A prefix probability becomes a "strong prefix" if the same

TABLE 1
A. ADMISSIBLE INITIAL CONSONANT STRINGS OF CVC WORDS

B	N	BL	GL	SH	TR	SCH
C	P	BR	GN	SK	TW	SCR
D	Q	CH	GR	SL	WH	SHR
F	R	CL	KN	SM	WR	SPH
G	S	CR	KR	SN		SPL
H	T	DR	PH	SP		SPR
J	V	DW	PL	SQ		STR
K	W	FL	PR	ST		THR
L	Z	FR	RH	SW		THW
M		GH	SC	TH		

B. ADMISSIBLE FINAL CONSONANT STRINGS OF CVC WORDS NOT ENDING WITH E

B	BB	LT	RF	WD	GHT
C	CH	MB	RK	WK	LCH
D	CK	MM	RL	WL	LPH
F	CT	MN	RM	WN	LTH
G	DD	MP	RN	XT	MPH
H	FF	ND	RP	ZZ	MPT
K	FT	NG	RR		NCH
L	GG	NK	RT		NTH
M	GH	NN	SH		NTZ
N	GN	NT	SK		RCH
P	LD	NX	SM		RSH
R	LF	PH	SP		RST
T	LK	PT	SS		RTH
W	LL	RB	ST		SCH
X	LM	RC	TH		TCH
Z	LP	RD	TT		

prefix probability arises from two or more inadmissible consonant strings. The definition for strong suffixes is analogous, proceeding from the other end of the word. Thus, given a word of the form $\dots V_3C_3V_2C_2V_1C_1$, if either C_2 or C_3 is an inadmissible string, there is a mandatory syllabic break within the string, and everything following that break is defined as a "suffix possibility." Then the definition for suffix probability and for strong suffix is the same as for prefixes above, in

* This work was accomplished under the Office of Naval Research and the Lockheed Independent Research Program. The author wishes to thank Dan L. Smith for writing many of the computer programs used in deriving the affixes.

¹ J. L. Dolby and H. L. Resnikoff, "The Nature of Affixing in Written English," *Mechanical Translation*, Vol. 8, Nos. 3, 4 (June and October, 1965), pp. 84-89.

which the word *suffix* can be substituted for the word *prefix* wherever it occurs. The consonant string C_1 may be blank in either case. The criterion of four or more words in establishing an affix probability and of two or more consonant strings in defining an affix from a probability was established by Dolby and Resnikoff. This criterion was established heuristically and has been retained here not only for the sake of consistency but also because it was proven effective.

The definition for weak affixes has also been extended to include two-syllable affixes. Weak affixes are so classified because their definition is based on a probable syllabic break rather than on a mandatory one. Because such probable breaks are not interior to a consonant string, weak prefixes end with a vowel and weak suffixes begin with one. For prefixes, given a word of the form $C_1V_1C_2V_2C_3V_3 \dots$, if either C_2 or C_3 is an admissible initial string but not an admissible final string, everything preceding that consonant string is a prefix possibility. For suffixes, given a word of the form $\dots V_3C_3V_2C_2V_1C_1$, if either C_2 or C_3 is an admissible final string but not an admissible initial string, everything following that consonant string is a suffix possibility. The criterion by which an affix possibility becomes an affix is the same as for strong affixes. Note that these definitions exclude admissible final strings from C_2 or C_3 for prefixes, and admissible initial strings from C_2 or C_3 for suffixes, in order to increase the reliability of the definition by reducing the probability of postulating a break before (for prefixes) or after (for suffixes) C_2 or C_3 where it does not exist. Consider the prefix case first. If C_2 or C_3 is an admissible initial string, and also an admissible ending string, the syllabic break could be logically either before or after the string. The string *CH* is such a string, as the following words illustrate:

enrich/ment	ta/chometer
poach/er	re/christen

By eliminating such doubtful strings we should increase somewhat the reliability of the definition of our prefix possibilities, but we do not completely eliminate chance for error, because even with initial strings not also final strings, a break may occur internal to a multi-letter string or after a single-letter string. The strings *BR* and *GR* are such multi-letter strings, as the following words illustrate:

sub/routine	ag/riculture
re/broadcast	de/gree

The chances of this happening in two multi-letter strings with the same prefix possibility is judged small enough to be discounted, since we are here simply defining prefix sequences. The chances of error due to a break after a single letter seems greater, as with the letter *S*:

re/sidual
res/ident

However, since there are only three single consonants that are beginning but not ending strings (*J*, *S*, *V*), and since again it takes two consonant strings to cause a sequence to be defined as an affix, this problem too can be discounted.

It is suspected that the situation for suffixes is more difficult in that the set of terminal consonant strings left after removing initial strings has more members that show a tendency to break internally. For example, breaks in the following strings are common:

<i>c/t</i> as in lac/tate	<i>m/b</i> as in am/bition
<i>r/t</i> as in fer/tile	<i>m/p</i> as in am/perere
<i>p/t</i> as in ap/titude	<i>r/l</i> as in pur/loin
<i>r/b</i> as in ar/bor	<i>n/d</i> as in ban/dit

and so on. Therefore, more difficulty in determining when a defined weak suffix is actually acting as a suffix in a given word could reasonably be anticipated. It would be interesting to subject each of the weak suffixes to a qualifying test, namely, that in the two-syllable data set there not be two sets of illegal strings preceding the suffix, where each set had at least four members. When this test was applied to the five suffixes *a*, *age*, *ah*, *ent*, and *ock*, two of the suffixes, *a* and *ock*, failed the test. But both *a* and *ock* obviously sometimes act as suffixes (they are both listed in the dictionaries as such), so it is unwise to eliminate them at this point in the research. What is indicated, perhaps, is the structural classification of the weak suffixes by degree of weakness as a means of approaching the suffix-in-context problem.

Table 2, which reviews the prefixes and suffixes defined by Resnikoff and Dolby, uses the two-vowel-string words as the data set. Table 3 shows the new suffixes defined using four-, five-, six-, and seven-vowel-string words, with the preceding letter strings and occurrence counts that established them as suffixes. Surprisingly, there is only one that can be considered a strong suf-

TABLE 2
AFFIXES FROM TWO-VOWEL-STRING WORDS

PREFIXES		SUFFIXES			
Strong	Weak	Strong	Weak		
ac	a	ful	a	eon	ive
ad	be	land	age	er	o
al	cy	ler	ah	et	ock
con	de	less	al	ic	on
dis	e	let	an	ie	or
en	i	ling	ant	ier	ot
ex	re	lock	ar	ile	ow
in		ly	ard	in	ue
mis		man	at	ine	um
out		ment	ed	ing	ure
sub		ness	ee	ion	us
sun		ward	el	is	
trans			en	ish	
us			ent	ite	

TABLE 3
SUFFIXES FROM MULTISYLLABLE WORDS

Suffix	Preceding Letter String	No. of Occurrences of Suffix Following the Given Letter String
(t)ation	<i>c(t)</i>	5
	<i>l(t)</i>	6
	<i>n(t)</i>	36
	<i>r(t)</i>	5
able	<i>ll</i>	8
	<i>nt</i>	4
ial	<i>nn</i>	8
	<i>nt</i>	37
ate	<i>ll</i>	6
	<i>nn</i>	5
ist	<i>ll</i>	4
	<i>nt</i>	12
	<i>pt</i>	4
ism	<i>ll</i>	4
	<i>nt</i>	5
ian	<i>ll</i>	4
	<i>nn</i>	4
ium	<i>ng</i>	5
	<i>rd</i>	4
ia	<i>ps</i>	12
	<i>rd</i>	4
	<i>nt</i>	5
y	<i>rg</i>	4
	<i>ps</i>	4
	<i>rm</i>	4
	<i>rr</i>	36
	<i>st</i>	19
ous	<i>x</i>	13
	<i>ll</i>	6
ide	<i>rm</i>	6
	<i>rp</i>	11
	<i>x</i>	9
is	<i>lf</i>	7
	<i>ps</i>	6
	<i>x</i>	28

fix, and that actually turned up as the weak suffix *ation*. Since all of the preceding letter strings turned out to be of the form *Ct* (where *C* = *c*, *l*, *n*, or *r*), and since phonetic breaks were consistently before the *t* (as in *plantation*), it seemed reasonable to consider *tation* a strong suffix. Of the thirteen newly defined suffixes, *able*, *ial*, *ate*, *ist*, *ism*, *y*, *ous*, *ian*, *ium*, *ia*, and *ide* are all commonly recognized as such, while only *tation* or *ation* and *is* are not.

It was expected that more than one two-vowel-string suffix would be obtained. Instead, a number of sequences were observed that appear to act as inner suffixes, or suffix-compounding elements, which occur frequently in combination with one-syllable suffixes. Thus, the sequence *tic* is frequently encountered followed by *al*, *ize*, or *ide* to form *tical*, *ticism*, *ticize*, or *ticide*, as in

TABLE 4
ELEMENTS COMBINING WITH SUFFIXES

Suffix-Compounding Element	Associated Terminal Letter String	No. of Occurrences
-cat-	<i>rc</i>	9
	<i>nc</i>	12
-mat-	<i>rm</i>	22
	<i>mm</i>	18
-pos-	<i>mp</i>	8
	<i>rp</i>	6
-pat-	<i>ip</i>	6
	<i>rp</i>	6
-sit-	<i>ns</i>	8
	<i>rs</i>	5
	<i>ss</i>	5
-sat-	<i>ns</i>	12
	<i>rs</i>	5
	<i>ss</i>	16
-tat-	<i>lt</i>	16
	<i>nt</i>	46
	<i>rt</i>	11
-tur-	<i>ct</i>	6
	<i>et</i>	19
	<i>nt</i>	8
-tic-	<i>ct</i>	13
	<i>nt</i>	7
	<i>pt</i>	13
-tor-	<i>ct</i>	33
	<i>nt</i>	6
-ter-	<i>ct</i>	8
	<i>et</i>	9
	<i>nt</i>	8
-tin-	<i>pt</i>	44
	<i>nt</i>	6
	<i>rt</i>	6

elliptical, *asepticism*, *didacticism*, *ascepticize*, *romanticize*, and *infanticide*. Such interior sequences that meet the occurrence criteria set up for suffixes are listed in Table 4. It is expected that these sequences will have little syntactic meaning but may be helpful in word-hyphenation techniques.

Table 5 shows the prefixes defined using four-, five-, six-, and seven-vowel-string words, with the following letter strings and occurrence counts that established them as prefixes. The three newly defined strong two-syllable prefixes *circum*, *inter*, and *hyper*, are well known. Three other common prefixes, *over*, *under*, and *super*, were encountered with a good many letter strings but always failed to meet the requirement of more than three occurrences with a given letter string.

Of the strong one-syllable prefixes defined, *ab*, *at*, *ap*, *com*, *an*, *em*, *im*, and *ec* are recognized by dictionaries, while *vul* is not. Of the weak two-syllable prefixes, *auto*, *demo*, *iso*, *photo*, *epi*, and *tele* are com-

monly recognized, but *ana*, *apo*, *deni*, and *irre* are not. (*Irre* is no doubt a combination of the recognized prefixes *i* and *re*.) None of the one-syllable weak prefixes (*au*, *ca*, *hy*, *ma*, *mi*, *lu*, *pro*, *sa*, *su*, *vi*) is familiar as a meaningful prefix except for *pro*. Therefore, the next

TABLE 5
PREFIXES FROM MULTISYLLABLE WORDS
A. WEAK PREFIXES

Prefix	Following Letter String	No. of Occurrences of Prefix Preceding the Given Letter String
ana	{ <i>cl</i>	4
	{ <i>gl</i>	6
apo	{ <i>cr</i>	4
	{ <i>str</i>	4
auto	{ <i>cr</i>	4
	{ <i>gr</i>	4
	{ <i>tr</i>	4
deni	{ <i>gr</i>	4
	{ <i>tr</i>	8
demo	{ <i>cr</i>	12
	{ <i>gr</i>	6
epi	{ <i>gr</i>	12
	{ <i>sc</i>	7
	{ <i>cl</i>	4
irre	{ <i>fr</i>	5
	{ <i>pr</i>	7
	{ <i>tr</i>	6
iso	{ <i>cr</i>	4
	{ <i>tr</i>	4
photo	{ <i>gr</i>	7
	{ <i>tr</i>	4
tele	{ <i>gr</i>	8
	{ <i>sc</i>	6
au	{ <i>sc</i>	6
	{ <i>str</i>	5
ca	{ <i>j</i>	4
	{ <i>pr</i>	4
	{ <i>sc</i>	5
hy	{ <i>dr</i>	85
	{ <i>gl</i>	5
ma	{ <i>cr</i>	18
	{ <i>j</i>	9
	{ <i>tr</i>	15
mi	{ <i>cr</i>	69
	{ <i>thr</i>	4
pro	{ <i>gl</i>	6
	{ <i>pr</i>	4
sa	{ <i>cr</i>	8
	{ <i>pr</i>	5
su	{ <i>bl</i>	6
	{ <i>pr</i>	11
	{ <i>sc</i>	5
vi	{ <i>br</i>	8
	{ <i>sc</i>	5
	{ <i>tr</i>	4

B. STRONG PREFIXES

Prefix	Defining Letter String	No. of Occurrences of Prefix with Given Letter String
at	{ <i>tm</i>	15
	{ <i>ttr</i>	11
ap	{ <i>ppl</i>	15
	{ <i>ppr</i>	46
an	{ <i>ndr</i>	18
	{ <i>ngl</i>	9
	{ <i>nh</i>	6
	{ <i>nih</i>	20
	{ <i>nthr</i>	35
em	{ <i>mbl</i>	12
	{ <i>mbr</i>	20
im	{ <i>mbr</i>	5
	{ <i>mpl</i>	21
	{ <i>mpr</i>	66
com	{ <i>mpl</i>	28
	{ <i>mpr</i>	13
vul	{ <i>lc</i>	6
	{ <i>ln</i>	4
ec	{ <i>cc</i>	4
	{ <i>ccl</i>	10
	{ <i>cst</i>	4
ob	{ <i>bj</i>	19
	{ <i>bs</i>	21
	{ <i>bsc</i>	9
	{ <i>bst</i>	6
	{ <i>bstr</i>	5
ab	{ <i>bt</i>	6
	{ <i>bd</i>	4
ab	{ <i>bn</i>	7
	{ <i>bs</i>	19
	{ <i>mc</i>	5
circum	{ <i>mf</i>	7
	{ <i>mr</i>	5
	{ <i>mscr</i>	6
	{ <i>mst</i>	10
	{ <i>mv</i>	10
inter	{ <i>rcl</i>	5
	{ <i>rj</i>	6
	{ <i>rpr</i>	9
	{ <i>rsp</i>	6
hyper	{ <i>rcr</i>	4
	{ <i>rpl</i>	4
	{ <i>rtr</i>	5

step, in which the part of speech implications of the structurally defined affixes is investigated, will be especially interesting for this group. It is, in fact, in the next steps, in which the various applications and implications of the structurally defined affixes are investigated, that the utility, and therefore the validity, of these structural definitions will be tested.

Received December 8, 1965