

# Writing Tools - The STYLE and DICTION Programs

*L. L. Cherry*

*W. Vesterman*

Livingston College  
Rutgers University

## ABSTRACT

Text processing systems are now in heavy use in many companies to format documents. With many documents stored on line, it has become possible to use computers to study writing style itself and to help writers produce better written and more readable prose. The system of programs described here is an initial step toward such help. It includes programs and a data base designed to produce a stylistic profile of writing at the word and sentence level. The system measures readability, sentence and word length, sentence type, word usage, and sentence openers. It also locates common examples of wordy phrasing and bad diction. The system is useful for evaluating a document's style, locating sentences that may be difficult to read or excessively wordy, and determining a particular writer's style over several documents.

## 1. Introduction

Computers have become important in the document preparation process, with programs to check for spelling errors and to format documents. As the amount of text stored on line increases, it becomes feasible and attractive to study writing style and to attempt to help the writer in producing readable documents. The system of writing tools described here is a first step toward such help. The system includes programs and a data base to analyze writing style at the word and sentence level. We use the term "style" in this paper to describe the results of a writer's particular choices among individual words and sentence forms. Although many judgements of style are subjective, particularly those of word choice, there are some objective measures that experts agree lead to good style. Three programs have been written to measure some of the objectively definable characteristics of writing style and to identify some commonly misused or unnecessary phrases. Although a document that conforms to the stylistic rules is not guaranteed to be coherent and readable, one that violates all of the rules is likely to be difficult or tedious to read. The program STYLE calculates readability, sentence length variability, sentence type, word usage and sentence openers at a rate of about 400 words per second on a PDP11/70 running the UNIX® Operating System. It assumes that the sentences are well-formed, i. e. that each sentence has a verb and that the subject and verb agree in number. DICTION identifies phrases that are either bad usage or unnecessarily wordy. EXPLAIN acts as a thesaurus for the phrases found by DICTION. Sections 2, 3, and 4 describe the programs; Section 5 gives the results on a cross-section of technical documents; Section 6 discusses accuracy and problems; Section 7 gives implementation details.

## 2. STYLE

The program STYLE reads a document and prints a summary of readability indices, sentence length and type, word usage, and sentence openers. It may also be used to locate all sentences in a document longer than a given length, of readability index higher than a given number, those containing a passive verb, or those beginning with an expletive. STYLE is based on the system for finding English word classes or parts of speech, PARTS [1]. PARTS is a set of programs that uses a small dictionary (about 350 words) and suffix rules to partially assign word classes to English text. It then uses experimentally derived rules of word order to assign word classes to all words in the text with an accuracy of about 95%. Because PARTS uses only a small dictionary and general rules, it works on text about any subject, from physics to psychology. Style measures have been built into the output phase of the programs that make up PARTS. Some of the measures are simple counters of the word classes found by PARTS; many are more complicated. For example, the verb count is the total number of verb phrases. This includes phrases like:

has been going  
was only going  
to go

each of which each counts as one verb. Figure 1 shows the output of STYLE run on a paper by Kernighan and Mashey about the UNIX programming environment [2].

programming environment	
readability grades:	
	(Kincaid) 12.3 (auto) 12.8 (Coleman-Liau) 11.8 (Flesch) 13.5 (46.3)
sentence info:	
	no. sent 335 no. wds 7419
	av sent leng 22.1 av word leng 4.91
	no. questions 0 no. imperatives 0
	no. nonfunc wds 4362 58.8% av leng 6.38
	short sent (<17) 35% (118) long sent (>32) 16% (55)
	longest sent 82 wds at sent 174; shortest sent 1 wds at sent 117
sentence types:	
	simple 34% (114) complex 32% (108)
	compound 12% (41) compound-complex 21% (72)
word usage:	
	verb types as % of total verbs
	tobe 45% (373) aux 16% (133) inf 14% (114)
	passives as % of non-inf verbs 20% (144)
	types as % of total
	prep 10.8% (804) conj 3.5% (262) adv 4.8% (354)
	noun 26.7% (1983) adj 18.7% (1388) pron 5.3% (393)
	nominalizations 2 % (155)
sentence beginnings:	
	subject opener: noun (63) pron (43) pos (0) adj (58) art (62) tot 67%
	prep 12% (39) adv 9% (31)
	verb 0% (1) sub_conj 6% (20) conj 1% (5)
	expletives 4% (13)

Figure 1

As the example shows, STYLE output is in five parts. After a brief discussion of sentences, we will describe the parts in order.

## 2.1. What is a sentence?

Readers of documents have little trouble deciding where the sentences end. People don't even have to stop and think about uses of the character "." in constructions like 1.25, A. J. Jones, Ph.D., i. e., or etc. . When a computer reads a document, finding the end of sentences is not as easy. First we must throw away the printer's marks and formatting commands that litter the text in computer form. Then STYLE defines a sentence as a string of words ending in one of:

. ! ? /.

The end marker "/" may be used to indicate an imperative sentence. Imperative sentences that are not so marked are not identified as imperative. STYLE properly handles numbers with embedded decimal points and commas, strings of letters and numbers with embedded decimal points used for naming computer file names, and the common abbreviations listed in Appendix 1. Numbers that end sentences, like the preceding sentence, cause a sentence break if the next word begins with a capital letter. Initials only cause a sentence break if the next word begins with a capital and is found in the dictionary of function words used by PARTS. So the string

J. D. JONES

does not cause a break, but the string

... system H. The ...

does. With these rules most sentences are broken at the proper place, although occasionally either two sentences are called one or a fragment is called a sentence. More on this later.

## 2.2. Readability Grades

The first section of STYLE output consists of four readability indices. As Klare points out in [3] readability indices may be used to estimate the reading skills needed by the reader to understand a document. The readability indices reported by STYLE are based on measures of sentence and word lengths. Although the indices may not measure whether the document is coherent and well organized, experience has shown that high indices seem to be indicators of stylistic difficulty. Documents with short sentences and short words have low scores; those with long sentences and many polysyllabic words have high scores. The 4 formulae reported are Kincaid Formula [4], Automated Readability Index [5], Coleman-Liau Formula [6] and a normalized version of Flesch Reading Ease Score [7]. The formulae differ because they were experimentally derived using different texts and subject groups. We will discuss each of the formulae briefly; for a more detailed discussion the reader should see [3].

The Kincaid Formula, given by:

$$Reading\_Grade = 11.8 * syl\_per\_wd + .39 * wds\_per\_sent - 15.59$$

was based on Navy training manuals that ranged in difficulty from 5.5 to 16.3 in reading grade level. The score reported by this formula tends to be in the mid-range of the 4 scores. Because it is based on adult training manuals rather than school book text, this formula is probably the best one to apply to technical documents.

The Automated Readability Index (ARI), based on text from grades 0 to 7, was derived to be easy to automate. The formula is:

$$Reading\_Grade = 4.71 * let\_per\_wd + .5 * wds\_per\_sent - 21.43$$

ARI tends to produce scores that are higher than Kincaid and Coleman-Liau but are usually slightly lower than Flesch.

The Coleman-Liau Formula, based on text ranging in difficulty from .4 to 16.3, is:

$$Reading\_Grade = 5.89 * let\_per\_wd - .3 * sent\_per\_100\_wds - 15.8$$

Of the four formulae this one usually gives the lowest grade when applied to technical

documents.

The last formula, the Flesch Reading Ease Score, is based on grade school text covering grades 3 to 12. The formula, given by:

$$\text{Reading\_Score} = 206.835 - 84.6 * \text{syl\_per\_wd} - 1.015 * \text{wds\_per\_sent}$$

is usually reported in the range 0 (very difficult) to 100 (very easy). The score reported by STYLE is scaled to be comparable to the other formulas, except that the maximum grade level reported is set to 17. The Flesch score is usually the highest of the 4 scores on technical documents.

Coke [8] found that the Kincaid Formula is probably the best predictor for technical documents; both ARI and Flesch tend to overestimate the difficulty; Coleman-Liau tend to underestimate. On text in the range of grades 7 to 9 the four formulas tend to be about the same. On easy text the Coleman-Liau formula is probably preferred since it is reasonably accurate at the lower grades and it is safer to present text that is a little too easy than a little too hard.

If a document has particularly difficult technical content, especially if it includes a lot of mathematics, it is probably best to make the text very easy to read, i.e. a lower readability index by shortening the sentences and words. This will allow the reader to concentrate on the technical content and not the long sentences. The user should remember that these indices are estimators; they should not be taken as absolute numbers. STYLE called with “-r number” will print all sentences with an Automated Readability Index equal to or greater than “number”.

### 2.3. Sentence length and structure

The next two sections of STYLE output deal with sentence length and structure. Almost all books on writing style or effective writing emphasize the importance of variety in sentence length and structure for good writing. Ewing’s first rule in discussing style in the book *Writing for Results* [9] is:

“Vary the sentence structure and length of your sentences.”

Leggett, Mead and Charvat break this rule into 3 in *Prentice-Hall Handbook for Writers* [10] as follows:

“34a. Avoid the overuse of short simple sentences.”

“34b. Avoid the overuse of long compound sentences.”

“34c. Use various sentence structures to avoid monotony and increase effectiveness.”

Although experts agree that these rules are important, not all writers follow them. Sample technical documents have been found with almost no sentence length or type variability. One document had 90% of its sentences about the same length as the average; another was made up almost entirely of simple sentences (80%).

The output sections labeled “sentence info” and “sentence types” give both length and structure measures. STYLE reports on the number and average length of both sentences and words, and number of questions and imperative sentences (those ending in “/.”). The measures of non-function words are an attempt to look at the content words in the document. In English non-function words are nouns, adjectives, adverbs, and non-auxiliary verbs; function words are prepositions, conjunctions, articles, and auxiliary verbs. Since most function words are short, they tend to lower the average word length. The average length of non-function words may be a more useful measure for comparing word choice of different writers than the total average word length. The percentages of short and long sentences measure sentence length variability. Short sentences are those at least 5 words less than the average; long sentences are those at least 10 words longer than the average. Last in the sentence information section is the length and location of the longest and shortest sentences. If the flag “-l number” is used, STYLE will print all sentences longer than “number”.

Because of the difficulties in dealing with the many uses of commas and conjunctions in English, sentence type definitions vary slightly from those of standard textbooks, but still

measure the same constructional activity.

1. A simple sentence has one verb and no dependent clause.
2. A complex sentence has one independent clause and one dependent clause, each with one verb. Complex sentences are found by identifying sentences that contain either a subordinate conjunction or a clause beginning with words like “that” or “who”. The preceding sentence has such a clause.
3. A compound sentence has more than one verb and no dependent clause. Sentences joined by “;” are also counted as compound.
4. A compound-complex sentence has either several dependent clauses or one dependent clause and a compound verb in either the dependent or independent clause.

Even using these broader definitions, simple sentences dominate many of the technical documents that have been tested, but the example in Figure 1 shows variety in both sentence structure and sentence length.

## 2.4. Word Usage

The word usage measures are an attempt to identify some other constructional features of writing style. There are many different ways in English to say the same thing. The constructions differ from one another in the form of the words used. The following sentences all convey approximately the same meaning but differ in word usage:

The cxio program is used to perform all communication between the systems.

The cxio program performs all communications between the systems.

The cxio program is used to communicate between the systems.

The cxio program communicates between the systems.

All communication between the systems is performed by the cxio program.

The distribution of the parts of speech and verb constructions helps identify overuse of particular constructions. Although the measures used by STYLE are crude, they do point out problem areas. For each category, STYLE reports a percentage and a raw count. In addition to looking at the percentage, the user may find it useful to compare the raw count with the number of sentences. If, for example, the number of infinitives is almost equal to the number of sentences, then many of the sentences in the document are constructed like the first and third in the preceding example. The user may want to transform some of these sentences into another form. Some of the implications of the word usage measures are discussed below.

### *Verbs*

are measured in several different ways to try to determine what types of verb constructions are most frequent in the document. Technical writing tends to contain many passive verb constructions and other usage of the verb “to be”. The category of verbs labeled “tobe” measures both passives and sentences of the form:

*subject tobe predicate*

In counting verbs, whole verb phrases are counted as one verb. Verb phrases containing auxiliary verbs are counted in the category “aux”. The verb phrases counted here are those whose tense is not simple present or simple past. It might eventually be useful to do more detailed measures of verb tense or mood. Infinitives are listed as “inf”. The percentages reported for these three categories are based on the total number of verb phrases found. These categories are not mutually exclusive; they cannot be added, since, for example, “to be going” counts as both “tobe” and “inf”. Use of these three types of verb constructions varies significantly among authors.

STYLE reports passive verbs as a percentage of the finite verbs in the document. Most style books warn against the overuse of passive verbs. Coleman [11] has shown that sentences with active verbs are easier to learn than those with passive verbs. Although the inverted object-subject order of the passive voice seems to emphasize the object, Coleman’s

experiments showed that there is little difference in retention by word position. He also showed that the direct object of an active verb is retained better than the subject of a passive verb. These experiments support the advice of the style books suggesting that writers should try to use active verbs wherever possible. The flag “-p” causes STYLE to print all sentences containing passive verbs.

#### *Pronouns*

add cohesiveness and connectivity to a document by providing back-reference. They are often a short-hand notation for something previously mentioned, and therefore connect the sentence containing the pronoun with the word to which the pronoun refers. Although there are other mechanisms for such connections, documents with no pronouns tend to be wordy and to have little connectivity.

#### *Adverbs*

can provide transition between sentences and order in time and space. In performing these functions, adverbs, like pronouns, provide connectivity and cohesiveness.

#### *Conjunctions*

provide parallelism in a document by connecting two or more equal units. These units may be whole sentences, verb phrases, nouns, adjectives, or prepositional phrases. The compound and compound-complex sentences reported under sentence type are parallel structures. Other uses of parallel structures are indicated by the degree that the number of conjunctions reported under word usage exceeds the compound sentence measures.

#### *Nouns and Adjectives.*

A ratio of nouns to adjectives near unity may indicate the over-use of modifiers. Some technical writers qualify every noun with one or more adjectives. Qualifiers in phrases like “simple linear single-link network model” often lend more obscurity than precision to a text.

#### *Nominalizations*

are verbs that are changed to nouns by adding one of the suffixes “ment”, “ance”, “ence”, or “ion”. Examples are accomplishment, admittance, adherence, and abbreviation. When a writer transforms a nominalized sentence to a non-nominalized sentence, she/he increases the effectiveness of the sentence in several ways. The noun becomes an active verb and frequently one complicated clause becomes two shorter clauses. For example,

Their inclusion of this provision is admission of the importance of the system.

When they included this provision, they admitted the importance of the system.

Coleman found that the transformed sentences were easier to learn, even when the transformation produced sentences that were slightly longer, provided the transformation broke one clause into two. Writers who find their document contains many nominalizations may want to transform some of the sentences to use active verbs.

## **2.5. Sentence openers**

Another agreed upon principle of style is variety in sentence openers. Because STYLE determines the type of sentence opener by looking at the part of speech of the first word in the sentence, the sentences counted under the heading “subject opener” may not all really begin with the subject. However, a large percentage of sentences in this category still indicates lack of variety in sentence openers. Other sentence opener measures help the user determine if there are transitions between sentences and where the subordination occurs. Adverbs and conjunctions at the beginning of sentences are mechanisms for transition between sentences. A pronoun at the beginning shows a link to something previously mentioned and indicates connectivity.

The location of subordination can be determined by comparing the number of sentences that begin with a subordinator with the number of sentences with complex clauses. If few sentences start with subordinate conjunctions then the subordination is embedded or at the end of the complex sentences. For variety the writer may want to transform some sentences to have leading subordination.

The last category of openers, expletives, is commonly overworked in technical writing. Expletives are the words “it” and “there”, usually with the verb “to be”, in constructions where the subject follows the verb. For example,

There are three streets used by the traffic.  
There are too many users on this system.

This construction tends to emphasize the object rather than the subject of the sentence. The flag “-e” will cause STYLE to print all sentences that begin with an expletive.

### 3. DICTION

The program DICTION prints all sentences in a document containing phrases that are either frequently misused or indicate wordiness. The program, an extension of Aho’s FGREP [12] string matching program, takes as input a file of phrases or patterns to be matched and a file of text to be searched. A data base of about 450 phrases has been compiled as a default pattern file for DICTION. Before attempting to locate phrases, the program maps upper case letters to lower case and substitutes blanks for punctuation. Sentence boundaries were deemed less critical in DICTION than in STYLE, so abbreviations and other uses of the character “.” are not treated specially. DICTION brackets all pattern matches in a sentence with the characters “[” “]” . Although many of the phrases in the default data base are correct in some contexts, in others they indicate wordiness. Some examples of the phrases and suggested alternatives are:

Phrase	Alternative
a large number of	many
arrive at a decision	decide
collect together	collect
for this reason	so
pertaining to	about
through the use of	by or with
utilize	use
with the exception of	except

Appendix 2 contains a complete list of the default file. Some of the entries are short forms of problem phrases. For example, the phrase “the fact” is found in all of the following and is sufficient to point out the wordiness to the user:

Phrase	Alternative
accounted for by the fact that	caused by
an example of this is the fact that	thus
based on the fact that	because
despite the fact that	although
due to the fact that	because
in light of the fact that	because
in view of the fact that	since
notwithstanding the fact that	although

Entries in Appendix 2 preceded by “~” are not matched. See Section 7 for details on the use of “~”.

The user may supply her/his own pattern file with the flag “-f patfile”. In this case the default file will be loaded first, followed by the user file. This mechanism allows users to suppress patterns contained in the default file or to include their own pet peeves that are not in the default file. The flag “-n” will exclude the default file altogether. In constructing a pattern file, blanks should be used before and after each phrase to avoid matching substrings in words. For example, to find all occurrences of the word “the”, the pattern “ the ” should be used. The blanks cause only the word “the” to be matched and not the string “the” in words like there,

other, and therefore. One side effect of surrounding the words with blanks is that when two phrases occur without intervening words, only the first will be matched.

#### 4. EXPLAIN

The last program, EXPLAIN, is an interactive thesaurus for phrases found by DICTION. The user types one of the phrases bracketed by DICTION and EXPLAIN responds with suggested substitutions for the phrase that will improve the diction of the document.

Table 1  
Text Statistics on 20 Technical Documents

	variable	minimum	maximum	mean	standard deviation
Readability	Kincaid	9.5	16.9	13.3	2.2
	automated	9.0	17.4	13.3	2.5
	Cole-Liau	10.0	16.0	12.7	1.8
	Flesch	8.9	17.0	14.4	2.2
sentence info.	av sent length	15.5	30.3	21.6	4.0
	av word length	4.61	5.63	5.08	.29
	av nonfunction length	5.72	7.30	6.52	.45
	short sent	23%	46%	33%	5.9
	long sent	7%	20%	14%	2.9
sentence types	simple	31%	71%	49%	11.4
	complex	19%	50%	33%	8.3
	compound	2%	14%	7%	3.3
	compound-complex	2%	19%	10%	4.8
verb types	tobe	26%	64%	44.7%	10.3
	auxiliary	10%	40%	21%	8.7
	infinitives	8%	24%	15.1%	4.8
	passives	12%	50%	29%	9.3
word usage	prepositions	10.1%	15.0%	12.3%	1.6
	conjunction	1.8%	4.8%	3.4%	.9
	adverbs	1.2%	5.0%	3.4%	1.0
	nouns	23.6%	31.6%	27.8%	1.7
	adjectives	15.4%	27.1%	21.1%	3.4
	pronouns	1.2%	8.4%	2.5%	1.1
	nominalizations	2%	5%	3.3%	.8
sentence openers	prepositions	6%	19%	12%	3.4
	adverbs	0%	20%	9%	4.6
	subject	56%	85%	70%	8.0
	verbs	0%	4%	1%	1.0
	subordinating conj	1%	12%	5%	2.7
	conjunctions	0%	4%	0%	1.5
	expletives	0%	6%	2%	1.7

#### 5. Results

##### 5.1. STYLE

To get baseline statistics and check the program's accuracy, we ran STYLE on 20 technical documents. There were a total of 3287 sentences in the sample. The shortest document was 67

sentences long; the longest 339 sentences. The documents covered a wide range of subject matter, including theoretical computing, physics, psychology, engineering, and affirmative action. Table 1 gives the range, median, and standard deviation of the various style measures. As you will note most of the measurements have a fairly wide range of values across the sample documents.

As a comparison, Table 2 gives the median results for two different technical authors, a sample of instructional material, and a sample of the Federalist Papers. The two authors show similar styles, although author 2 uses somewhat shorter sentences and longer words than author 1. Author 1 uses all types of sentences, while author 2 prefers simple and complex sentences, using few compound or compound-complex sentences. The other major difference in the styles of these authors is the location of subordination. Author 1 seems to prefer embedded or trailing subordination, while author 2 begins many sentences with the subordinate clause. The documents tested for both authors 1 and 2 were technical documents, written for a technical audience. The instructional documents, which are written for craftspeople, vary surprisingly little from the two technical samples. The sentences and words are a little longer, and they contain many passive and auxiliary verbs, few adverbs, and almost no pronouns. The instructional documents contain many imperative sentences, so there are many sentence with verb openers. The sample of Federalist Papers contrasts with the other samples in almost every way.

Table 2  
Text Statistics on Single Authors

	variable	author 1	author 2	inst.	FED
readability	Kincaid	11.0	10.3	10.8	16.3
	automated	11.0	10.3	11.9	17.8
	Coleman-Liau	9.3	10.1	10.2	12.3
	Flesch	10.3	10.7	10.1	15.0
sentence info	av sent length	22.64	19.61	22.78	31.85
	av word length	4.47	4.66	4.65	4.95
	av nonfunction length	5.64	5.92	6.04	6.87
	short sent	35%	43%	35%	40%
	long sent	18%	15%	16%	21%
sentence types	simple	36%	43%	40%	31%
	complex	34%	41%	37%	34%
	compound	13%	7%	4%	10%
	compound-complex	16%	8%	14%	25%
verb type	tobe	42%	43%	45%	37%
	auxiliary	17%	19%	32%	32%
	infinitives	17%	15%	12%	21%
	passives	20%	19%	36%	20%
word usage	prepositions	10.0%	10.8%	12.3%	15.9%
	conjunctions	3.2%	2.4%	3.9%	3.4%
	adverbs	5.05%	4.6%	3.5%	3.7%
	nouns	27.7%	26.5%	29.1%	24.9%
	adjectives	17.0%	19.0%	15.4%	12.4%
	pronouns	5.3%	4.3%	2.1%	6.5%
	nominalizations	1%	2%	2%	3%
sentence openers	prepositions	11%	14%	6%	5%
	adverbs	9%	9%	6%	4%
	subject	65%	59%	54%	66%
	verb	3%	2%	14%	2%
	subordinating conj	8%	14%	11%	3%
	conjunction	1%	0%	0%	3%
	expletives	3%	3%	0%	3%

## 5.2. DICTION

In the few weeks that DICTION has been available to users about 35,000 sentences have been run with about 5,000 string matches. The authors using the program seem to make the suggested changes about 50-75% of the time. To date, almost 200 of the 450 strings in the default file have been matched. Although most of these phrases are valid and correct in some contexts, the 50-75% change rate seems to show that the phrases are used much more often than concise diction warrants.

## 6. Accuracy

### 6.1. Sentence Identification

The correctness of the STYLE output on the 20 document sample was checked in detail. STYLE misidentified 129 sentence fragments as sentences and incorrectly joined two or more sentences 75 times in the 3287 sentence sample. The problems were usually because of non-standard formatting commands, unknown abbreviations, or lists of non-sentences. An impossibly long sentence found as the longest sentence in the document usually is the result of a long list of non-sentences.

### 6.2. Sentence Types

Style correctly identified sentence type on 86.5% of the sentences in the sample. The type distribution of the sentences was 52.5% simple, 29.9% complex, 8.5% compound and 9% compound-complex. The program reported 49.5% simple, 31.9% complex, 8% compound and 10.4% compound-complex. Looking at the errors on the individual documents, the number of simple sentences was under-reported by about 4% and the complex and compound-complex were over-reported by 3% and 2%, respectively. The following matrix shows the programs output vs. the actual sentence type.

		Program Results			
		simple	complex	compound	comp-complex
Actual Sentence Type	simple	1566	132	49	17
	complex	47	892	6	65
	compound	40	6	207	23
	comp-complex	0	52	5	249

The system's inability to find imperative sentences seems to have little effect on most of the style statistics. A document with half of its sentences imperative was run, with and without the imperative end marker. The results were identical except for the expected errors of not finding verbs as sentence openers, not counting the imperative sentences, and a slight difference (1%) in the number of nouns and adjectives reported.

### 6.3. Word Usage

The accuracy of identifying word types reflects that of PARTS, which is about 95% correct. The largest source of confusion is between nouns and adjectives. The verb counts were checked on about 20 sentences from each document and found to be about 98% correct.

## 7. Technical Details

### 7.1. Finding Sentences

The formatting commands embedded in the text increase the difficulty of finding sentences. Not all text in a document is in sentence form; there are headings, tables, equations and lists, for example. Headings like “Finding Sentences” above should be discarded, not attached to the next sentence. However, since many of the documents are formatted to be phototypeset, and contain font changes, which usually operate on the most important words in the document, discarding all formatting commands is not correct. To improve the programs’ ability to find sentence boundaries, the deformatting program, DEROFF [13], has been given some knowledge of the formatting packages used on the UNIX operating system. DEROFF will now do the following:

1. Suppress all formatting macros that are used for titles, headings, author’s name, etc.
2. Suppress the arguments to the macros for titles, headings, author’s name, etc.
3. Suppress displays, tables, footnotes and text that is centered or in no-fill mode.
4. Substitute a place holder for equations and check for hidden end markers. The place holder is necessary because many typists and authors use the equation setter to change fonts on important words. For this reason, header files containing the definition of the EQN delimiters must also be included as input to STYLE. End markers are often hidden when an equation ends a sentence and the period is typed inside the EQN delimiters.
5. Add a "." after lists. If the flag `-ml` is also used, all lists are suppressed. This is a separate flag because of the variety of ways the list macros are used. Often, lists are sentences that should be included in the analysis. The user must determine how lists are used in the document to be analyzed.

Both STYLE and DICTON call DEROFF before they look at the text. The user should supply the `-ml` flag if the document contains many lists of non-sentences that should be skipped.

### 7.2. Details of DICTON

The program DICTON is based on the string matching program FGREP. FGREP takes as input a file of patterns to be matched and a file to be searched and outputs each line that contains any of the patterns with no indication of which pattern was matched. The following changes have been added to FGREP:

1. The basic unit that DICTON operates on is a sentence rather than a line. Each sentence that contains one of the patterns is output.
2. Upper case letters are mapped to lower case.
3. Punctuation is replaced by blanks.
4. All pattern matches in the sentence are found and surrounded with “[” “]” .
5. A method for suppressing a string match has been added. Any pattern that begins with “~” will not be matched. Because the matching algorithm finds the longest substring, the suppression of a match allows words in some correct contexts not to be matched while allowing the word in another context to be found. For example, the word “which” is often incorrectly used instead of “that” in restrictive clauses. However, “which” is usually correct when preceded by a preposition or “,”. The default pattern file suppresses the match of the common prepositions or a double blank followed by “which” and therefore matches only the suspect uses. The double blank accounts for the replaced comma.

## **8. Conclusions**

A system of writing tools that measure some of the objective characteristics of writing style has been developed. The tools are sufficiently general that they may be applied to documents on any subject with equal accuracy. Although the measurements are only of the surface structure of the text, they do point out problem areas. In addition to helping writers produce better documents, these programs may be useful for studying the writing process and finding other formulae for measuring readability.

## References

1. L. L. Cherry, "PARTS - A System for Assigning Word Classes to English Text," submitted *Communications of the ACM*.
2. B. W. Kernighan and J. R. Mashey, "The UNIX Programming Environment," *Software – Practice & Experience*, **9**, 1-15 (1979).
3. G. R. Klare, "Assessing Readability," *Reading Research Quarterly*, 1974-1975, **10**, 62-102.
4. E. A. Smith and P. Kincaid, "Derivation and validation of the automated readability index for use with technical materials," *Human Factors*, 1970, **12**, 457-464.
5. J. P. Kincaid, R. P. Fishburne, R. L. Rogers, and B. S. Chissom, "Derivation of new readability formulas (Automated Readability Index, Fog count, and Flesch Reading Ease Formula) for Navy enlisted personnel," Navy Training Command Research Branch Report 8-75, Feb., 1975.
6. M. Coleman and T. L. Liau, "A Computer Readability Formula Designed for Machine Scoring," *Journal of Applied Psychology*, 1975, **60**, 283-284.
7. R. Flesch, "A New Readability Yardstick," *Journal of Applied Psychology*, 1948, **32**, 221-233.
8. E. U. Coke, private communication.
9. D. W. Ewing, *Writing for Results*, John Wiley & Sons, Inc., New York, N. Y. (1974).
10. G. Leggett, C. D. Mead and W. Charvat, *Prentice-Hall Handbook for Writers*, Seventh Edition, Prentice-Hall Inc., Englewood Cliffs, N. J. (1978).
11. E. B. Coleman, "Learning of Prose Written in Four Grammatical Transformations," *Journal of Applied Psychology*, 1965, vol. 49, no. 5, pp. 332-341.
12. A. V. Aho and M. J. Corasick, "Efficient String Matching: an aid to Bibliographic Search," *Communications of the ACM*, **18**, (6), 333-340, June 1975.
13. Bell Laboratories, "*UNIX TIME-SHARING SYSTEM: UNIX PROGRAMMER'S MANUAL*," Seventh Edition, Vol. 1 (January 1979).

Appendix 1

STYLE Abbreviations

a. d.  
A. M.  
a. m.  
b. c.  
Ch.  
ch.  
ckts.  
dB.  
Dept.  
dept.  
Depts.  
depts.  
Dr.  
Drs.  
e. g.  
Eq.  
eq.  
et al.  
etc.  
Fig.  
fig.  
Figs.  
figs.  
ft.  
i. e.  
in.  
Inc.  
Jr.  
jr.  
mi.  
Mr.  
Mrs.  
Ms.  
No.  
no.  
Nos.  
nos.  
P. M.  
p. m.  
Ph. D.  
Ph. d.  
Ref.  
ref.  
Refs.  
refs.  
St.  
vs.  
yr.

## Appendix 2

## Default DICTON Patterns

a great deal of	center portion	fearful that	in the form of
a large number of	check into	few in number	in the instance of
a lot of	check on	file away	in the interim
a majority of	check up on	final completion	in the last analysis
a need for	circle around	final ending	in the matter of
a number of	close proximity	final outcome	in the near future
a particular preference for	collaborate together	final result	in the neighborhood of
a preference for	collect together	finalize	in the not too distant future
a small number of	combine together	find it interesting to know	in the proximity of
a tendency to	come to an end	first and foremost	in the range of
abovementioned	commence	first beginnings	in the same way as described
absolutely complete	common accord	first initiated	in the shape of
absolutely essential	compensation	firstly	in the vicinity of
accomplished	completely eliminated	follow after	in this case
accordingly	comprise	following after	in view of the
activate	concerning	for the purpose of	in violation of
actual	conduct an investigation of	for the reason that	inasmuch as
added increments	conjecture	for the simple reason that	indicate
adequate enough	connect up	for this reason	indicative of
advent	consensus of opinion	for your information	initialize
afford an opportunity	consequent result	from the point of view of	initiate
aggregate	consolidate together	full and complete	injurious to
all of	construct	generally agreed	inquire
all throughout	contemplate	good and	inside of
along the line	continue on	got to	institute a
an indication of	continue to remain	gratuitous	intents and purposes
analyzation	could of	greatly minimize	intermingle
and etc	count up	head up	irregardless
and or	couple together	help but	is defined as
another additional	debate about	helps in the production of	is used to control
any and all	decide on	hopeful	is when
arrive at a	deleterious effect	if and when	is where
as a matter of fact	demean	if at all possible	it is incumbent
as a method of	demonstrate	impact	it stands to reason
as good or better than	depreciate in value	implement	it was noted that if
as of now	deserving of	important essentials	joint cooperation
as per	desirable benefits	importantly	joint partnership
as regards	desirous of	in a large measure	just exactly
as related to	different than	in a position to	kind of
as to	discontinue	in accordance	know about
assistance	disutility	in advance of	last but not least
assistance to	divide up	in agreement with	later on
assistance to	doubt but	in all cases	leaving out of consideration
assuming that	due to	in back of	liable
at a later date	duly noted	in behalf of	link up
at about	during the time that	in behind	literally
at above	each and every	in between	little doubt that
at all times	early beginnings	in case	lose out on
at an early date	effectuate	in close proximity	lots of
at below	emotional feelings	in conflict with	main essentials
at the present	empty out	in conjunction with	make a
at the time when	enclosed herein	in connection with	make adjustments to
at this point in time	enclosed herewith	in fact	make an
at this time	end result	in large measure	make application to
at which time	end up	in many cases	make contact with
at your earliest convenience	endeavor	in most cases	make mention of
authorization	enter in	in my opinion I think	make out a list of
awful	enter into	in order to	make the acquaintance of
basic fundamentals	enthused	in rare cases	make the adjustment
basically	entirely complete	in reference to	manner
be cognizant of	equally good as	in regard to	maximum possible
being as	essentially	in regards to	meaningful
being that	eventuate	in relation with	meet up with
brief in duration	every now and then	in short supply	melt down
bring to a conclusion	exactly identical	in size	melt up
but that	experiencing difficulty	in terms of	methodology
but what	fabricate	in the amount of	might of
by means of	face up to	in the case of	minimize as far as possible
by the use of	facilitate	in the course of	minor importance
carry out experiments	facts and figures	in the event	miss out on
center about	fast in action	in the field of	modification
center around	fearful of		

more preferable	seems apparent	worth while
most unique	send a communication	would of
must of	short space of time	ing behavior
mutual cooperation	should of	wise
necessary requisite	single unit	˘ which
necessitate	situation	˘ about which
need for	so as to	˘ after which
nice	sort of	˘ at which
not be un	spell out	˘ between which
not in a position to	still continue	˘ by which
not of a high order of accuracy	still remain	˘ for which
not un	subsequent	˘ from which
notwithstanding	substantially in agreement	˘ in which
of considerable magnitude	succeed in	˘ into which
of that	suggestive of	˘ of which
of the opinion that	superior than	˘ on which
off of	surrounding circumstances	˘ on which
on a few occasions	take appropriate	˘ over which
on account of	take cognizance of	˘ through which
on behalf of	take into consideration	˘ to which
on the grounds that	termed as	˘ under which
on the occasion	terminate	˘ upon which
on the part of	termination	˘ with which
one of the	the author	˘ without which
open up	the authors	˘clockwise
operates to correct	the case that	˘likewise
outside of	the fact	˘otherwise
over with	the foregoing	
overall	the foreseeable future	
past history	the fullest possible extent	
perceptive of	the majority of	
perform a measurement	the nature	
perform the measurement	the necessity of	
permits the reduction of	the only difference being that	
personalize	the order of	
pertaining to	the point that	
physical size	the truth is	
plan ahead	there are not many	
plan for the future	through the medium of	
plan in advance	through the use of	
plan on	throughout the entire	
present a conclusion	time interval	
present a report	to summarize the above	
presently	total effect of all this	
prior to	totality	
prioritize	transpire	
proceed to	true facts	
procure	try and	
productive of	ultimate end	
prolong the duration	under a separate cover	
protrude out from	under date of	
provided that	under separate cover	
pursuant to	under the necessity to	
put to use in	underlying purpose	
range all the way from	undertake a study	
reason is because	uniformly consistent	
reason why	unique	
recur again	until such time as	
reduce down	up to this time	
refer back	upshot	
reference to this	utilize	
reflective of	very	
regarding	very complete	
regretful	very unique	
reinitiate	vital	
relative to	which	
repeat again	with a view to	
representative of	with reference to	
resultant effect	with regard to	
resume again	with the exception of	
retreat back	with the object of	
return again	with the result that	
return back	with this in mind, it is clear that	
revert back	within the realm of possibility	
seal off	without further delay	