

## A TOOL THAT DETECTS PLAGIARISM IN PASCAL PROGRAMS

Sam Grier Department of Astronautics and Computer Science USAF ACADEMY USAFA CO 80840

Plagiarism has become a problem in introductory Computer Science courses. Programmed assignments can be copied and transformed with little human effort. A pertinent recommendation has resulted from this realization; an on-line system to detect programs that are "too similar" and hence suspected of plagiarism should be developed [4]. This paper discusses such a system for Pascal programs.

### 1. INTRODUCTION

As noted in recent literature, plagiarism has become a problem in introductory Computer Science courses [4]. To put it succinctly, students are copying other students' programs.

Detecting this plagriarism is difficult. Not only must graders grade a large volume of programs, but these programs all solve the same problem. Sophisticated plagiarism is not the problem; the sheer volume of code involved is simply overwhelming.

One attempted solution to this problem has been the development of a program at Purdue University by K.J. Ottenstein that quantifies the sameness of Fortran programs [3]. This program utilizes the four basic Software Science parameters suggested by M. Halstead as useful measures of program length [2]. This program utilizes only these parameters, and it counts them in a straightforward manner. The parameters are: the number of unique operators,
 the number of unique operands, (3) the total number of occurrences of operators, and (4) the total number of occurrences of operands [3]. It seems the first

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1981 ACM 0-89791-036-2/81/0200/0015 \$00.75

suggestion to use these parameters as measures of similarity or dissimilarity (depending on your viewpoint) came from N. Bulut as a by-product of his study of invariant properties of algorithms [1].

A tool that analyzes Pascal programs to detect those pairs of programs sufficiently similar such that plagiarism is a possibility is not known to exist. Program Accuse attempts to fill this void.

# 2. DESIGN

Program Accuse attempts to go beyond M. Halstead's four basic Software Science parameters in the belief that additional parameters are available to establish dissimilarity of two or more programs. It uses seven parameters and various counting heuristics that result in the computation of a correlation number that is used to determine the similarity of two programs. Accuse measures 20 parameters. The seven that comprise the correlation number were selected by testing different combinations of them.

An overriding concern of the development of Accuse has been that it be as inexpensive to use as possible. For this reason, the idea of utilizing the front end of a compiler was discarded, and Ottenstein's lead of using a fast counter was followed.

The result is a compromise between speed and comprehensive analysis. Accuse processes over 170 lines per second. However, it will not discover changes made by the sophisticated plagiarist. This is rationalized with the assumption that the student intelligent enough to plagiarize with sophistication has no need to plagiarize. We hope, however, that Accuse is not so simple-minded that it is easy to beat. It is meant to make plagiarism difficult to achieve, and it is meant to do this is such a manner that its repeated use does not compromise its heuristics. Accuse is a 2800 line program written in Pascal that runs on the E.T.H. Zuerich/ University of Minnesota Compiler. It consists of a modified Pascal scanner that passes tokens to a driver capable of processing compilable input programs. It also contains a host of support routines for the driver.

Accuse presently measures the following 20 parameters:

- 1. total lines
- 2. code lines
- 3. code comment lines
- 4. multiple statement lines
- 5. constants and types
- 6. variables declared (and used)
- 7. variables declared (and not used)
- 8. procedures and functions
- 9. var parameters
- 10. value parameters
- 11. procedure variables (includes 9 and 10)
- 12. for statements
- 13. repeat statements
- 14. while statements
- 15. goto statements
- 16. unique operators
- 17. unique operands
- 18. total operators
- 19. total operands
- 20. indenting function

The seven parameters that comprise the correlation number are:

- 1. unique operators
- 2. unique operands
- 3. total operators
- 4. total operands
- 5. code lines
- 6. variables declared (and used)

7. total control statements

One result of Accuse's development has been the failure of an "indenting function" to play a role in the detection of plagiarism. The indenting function is defined as:

- ( (left indentations) mod 1000) \*
  - 1000000 +
- ((right indentations) mod 1000) \* 1000 +
  - (zero indentations) mod 1000

If all programs were processed through a "pretty printer," an indenting function might become important. This additional cost is presently considered prohibitive, and it is contrary to the intent of Accuse being inexpensive to use.

The counting heuristics Accuse uses involve "total operators" and "code lines." "Total operators" does not include assignment operators. Additionally, for every assignment operator found, two operands are subtracted from "total operands," and "code lines" is decremented. This should prevent Accuse from being misled by unnecessary initializations and unnecessary assignment statements. "Code lines" ignores blank lines, comment lines, and declarations. It counts only executable lines of code within a program. "Code lines" was found to be an accurate indication of the sameness of two programs.

As Accuse only counts variables, the obvious tactic of changing variable names makes no difference to Accuse. Since Pascal requires declarations, Accuse can keep track of variables declared and subsequently used or not used. Hence excess declarations are an ineffective change to a program. Constants of enumerated types and tag fields in case clauses of record declarations that contain a declaration are considered variables. Since these constants cannot be read or written, their non-use is considered notable.

Accuse is also selective about what it calls operators. Software Science considers a BEGIN END combination as an operator [2]. Because BEGINs and ENDs can be added to Pascal code where not required, Accuse chooses to ignore them. Parentheses and several other operators are ignored by Accuse for essentially the same reason.

3. OUTPUT

Accuse prints four results for the user. The first (Table 1) is a dump of each program's identifier and its values of the 20 parameters measured by Accuse. This dump is sorted on the indenting function (a matter of my preference).

The second result (Table 2) is a

dump of each program's identifier and its respective values of the seven parameters used to compute the correlation number; each parameter list is sorted smallest to largest. In the output, the column headed FOR STMT actually contains the total number of control statements. This is a result of the implementation of summing parameters.

The third result (Table 4) is a frequency distribution graph that indicates the number of pairs of programs with like correlation numbers.

The final result (Table 5) is a list of all pairs of programs with correlation number greater than or equal to 28. Twenty nine is currently identified as the number that indicates the possibility of plagiarism, with 32 the maximum correlation number possible.

4. CORRELATION SCHEME

The scheme that computes the correlation number is only a tentative one. The current scheme was developed and tuned by using a group of 43 programs from an introductory course. Code for three of the programs was written together, but finished individually. The "importance values" for the seven correlation parameters were then adjusted until these three programs were brought into the domain of "those programs suspected of plagiarism."

The current correlation scheme involves computing an increment for each pair of affected programs based on the equation:

where pcounta and pcountb represent parameter counts, and (pcounta - pcountb) is less than or equal to some "window" size, depending on the particular parameter.

The computation of the correlation number may well be subject to improvement by a more elaborate scheme, or by simple changes to the importance values.

A simple, illustrative run of Accuse follows the text of this paper (Tables 1 through 5). This run processed 13 programs, three of which were input twice. Included is a print-out of the triangular matrix (Table 3) that contains correlation values of the pairs of programs. This matrix is not printed in a production model of Accuse.

Below we illustrate the computation of the correlation number for a pair of programs in the run. Before proceeding, it is necessary to note the following "window" sizes and "importance" values for each of the correlation parameters:

- 1. total operators
  window size = 5
  importance value = 6
- 2. total operands
  window size = 5
  importance value = 6
- 3. unique operators
  window size = 3
  importance value = 5
- 4. unique operands
  window size = 3
  importance value = 5
- 5. code lines window size = 3 importance value = 5
- 6. declared variables (and used)
  window size = 2
  importance value = 3
- 7. control statements
  window size = 1
  importance value = 2

The correlation number for the pair of programs T107 and T102 (see Tables) is computed as follows:

 T107 - T102 = 8 Eight is greater than the window size for this parameter, hence these are not "affected" programs.

 T107 - T102 = 16 Again, these are not "affected" programs.

3. T107 - T102 = 1 These programs are now within the window size, and an increment is calculated for this pair of programs. increment = 5 - (25 - 24) = 4

- correlation number = 4
- 4. T102 T107 = 0 increment = 5 - (13 - 13) = 5 correlation number = 9
- 5. T102 T107 = 1increment = 5 - (64 - 63) = 4 correlation number = 13
- 6. T107 T102 = 0 increment = 3 - (11 - 11) = 3 correlation number = 16
- 7. T102 T107 = 0 increment = 2 - (4 - 4) = 2 correlation number = 18

### 5. RESULTS

A typical production run of Accuse included 137 input programs consisting of 13,374 lines of code. Accuse processed the code on a CDC machine at a cost of \$12.32. It required:

> FL TO LOAD 110700 FL TO RUN 77100 105237B CM USED -89.956 CP SECS

Accuse prints all pairs of programs with correlation number greater than or equal to 28, though 29 is the number that indicates the possibility of plagiarism.

Several points are necessary.

In six runs of Accuse, sabotage occurred in two. There is nothing to prevent a student from removing lines of code from his program. One student shuffled his cards, and another added control characters not found in the character set of the machine. This led to Accuse being used in the context of a larger tool in its last and most successful run. The instructor retrieved the students' programs, compiled them, ran the programs on data the students had never seen, and then sent the source code to a file to be run on Accuse. This is the recommended context for the use of Accuse.

The correlation scheme is admittedly ad hoc. The only thing that can be said in its defense is that it seems to work. The use of Accuse should not be misunderstood. Accuse does not judge plagiarism; it merely indicates its possibility. It is a tool for the user to aid him in its detection; the decision as to plagiarism is left to the user. High correlation numbers may be meaningless; though rare, programs that are completely different may have like values for the seven parameters that are used to compute the correlation number.

6. THE REAL ISSUE?

Finally, as a reviewer noted, Accuse is a tool to discourage dishonesty in students. But, he asks, does anyone care to ask students why they cheat more now, and can we find ways to abort this rising phenomenon? These are pertinent educational issues.

7. ACKNOWLEDGEMENTS

.

Thanks to Lloyd Fosdick, who conceived this project and let me work on it; special thanks to Malcolm Newey for his insights and encouragement. 8. REFERENCES

1. Bulut, N., "Invariant Properties of Algorithms," PHD Thesis, Purdue University (August 1973) 118-119.

2. Halstead, M.H., "Elements of Software Science," Elsevier North Holland, New York (1977), Chapters 1-4,7.

3. Ottenstein, K.J., "An Algorithmic Approach to the Detection and Prevention of Plagiarism," SIGCSE Bulletin, Dec 76, Vol.8, No.4.

4. Shaw, M.; Jones, A; Knueven, P.; McDermott, J.; Miller, P.; Notkin, D., "Cheating Policy in a Computer Science Department," SIGCSE Bulletin, Jul 80, Vol. 12, No. 2.

\*\*\*\*\* Accuse was developed as a Masters Thesis under the advisement of Lloyd Fosdick, University of Colorado, Boulder.

SN 0		<b>;</b> ;	-	ī.	Ş	2	1997	С. Г	5	e M	24	5.4
ан <u>-</u>	• • •	Ś	2		•.		Ĵ		•	•	1	:
ī			~		••		<u>.</u>	·	•	÷,	٠.	*
101AL 0P1:P5	69 136	¥.	96	73	117	Ð	35	5 G 4	66	00	1 06	106
TOTAL OPERS	90 156	115	115	1 00	134	67	117	183	110	110	118	116
0170 0P105	4 1	13	13	12		13	ť:	13	5	13	1.3	<b></b>
58340 DINN	25 25	5	23	50	4	24	25	25	24	24	ម្នា. C+	25
GDTO SIMTS	00	0	0	0	0	0	0	0	0	0	à	0
WH 1 LE STOTS	ti ui	t	サ	1	4	ć	t	t•	4	4	1	4
REP STMTS	00	0	0	0	0	0	Q	0	0	0	Ċ	0
F.C.R Stwie	00	0	o	0	0	0	0	o	0	0	0	Ċ
PROC V43S	ເລ	on o	3	9	æ	ŝ	ŝ	æ	'n	80	8	n
LAU DARAN	<b>CV</b> (5)		(1	2	54	<	3	2	2	2	~	N
VAR Paraw	ო ო	• m	n	ę	e	m	m	(1)	m	e	'n	an.
PROCS AND FUNCS	c1 0	ŝ	~	C4	2	2	2	3	~	~	7	Ĩ
VARS NOT USED	00	0	0	o	0	0	0	0	0	0	0	ů
DECL VARS	10	-	:	01	:	:	:	<b>.</b>	Ξ	:	÷	
CONS AND TYPES	<b>ო</b> ო	ι <b>Μ</b>	e	ო	m	m	ო	m	ო	ო	e	ŝ
MULT STMT LINES	00	0	0	0	9	0	0	0	0	0	0	ç
CODE COMNT LINES	4 M 9 A 9 A	63	63	43	60	04	64	80	75	75	78	£8
CODE	48 10	58	58	71	66	44	56	19	64	64	63	63
TDTAL LINES	161 166	158	158	125	174	189	178	226	200	200	189	189
	1110	1106	1103	7112	1111	T1C9	1101	1108	T102	T105	T104	1107

.

VARIABLE(S) WERE DECLARED BUT NDT USED
 PROCEDURE(S) WERE DECLARED BUT NDT USED
 VARIABLE(S) AND PROCEDURE(S) WERE DECLARED BUT NDT USED

# TABLE 1

			.,	•3	• 1	4	17	7	et.	·1	F	មា	
S	1109	1+06	10.1	1113	1103	7104	7107	1102	105	T111	T 1 C B	1110	
DECL VARS	10	:	:	:	:		11			:	=	2	12
	1112	1162	1107	1105	1108	1111	T101	1109	11.6	T104	1103	<b>1113 156 1113 136 1101 25 1113 14 1113 72 1113 12 1110 5</b>	111.
COUE THES	44	18	56	58	58	5 <b>3</b>	63	64	64	66	71	72	۵, L
<b>!</b>	1:09	1110	1101	1106	1103	1104	1107	1102	1105	1111	T112	<b>T113 156</b> T113 136 T101 25 T113 14 T113 22 T113 12 T110 5	1108
CN10 CN02	12	13		13	175 175	<u>.</u>	13	13	13	13	t -	14	7
0	T112	T107	1108	1103	101	1106	1 109	1105	1104	1:02	1110	T113	1111
UNIQ FERS	24	24	24	24	54	2.E	25	, .	25	52	22	25	25
o	1102	1:05	T115	1112	1109	1101	T105	1113	1107	1106	1110	1:01	1103
DTAL	68	73	81	06	06	96	96	96	106	106	117	136	159
·- U	1110	1112	1109	1102	1105	1106	E017	1011	1107	T104	1111	1113	1108
DERS	06	57	100	110	110	115	115	117	118	118	134	156	183
μO	1110	1109	1112	1105	T102	1103	1106	T101	T104	1107	1111	T113	T106

WARNING: HEADING(S) ARE NOT CORRECT. AT LEAST ONE REFLECTS & SUM OF TWO COUNTS.

. TABLE 2													
	F W	n) +-	12	ŝ	÷	12	5	ŝ	ŝ		4	3	2
	F 0	0	-	10	10	Ξ	0	5	10	14	J	2	
•	****	13	17	÷	15		13	15	ü	÷	2		
	** 6	5		ç	ŭ	2	ç	2	(N 	0			
	F + O D		Ч ,	е Е	е Г	4	<u></u>	ĉ	<u>ن</u> 1				
	F-00	5	4	ŝ	ព	4	ŝ	5					
	+-01	50	0	<u>م</u>	32	с; т	18						
	H - 0 0	00	5	32	8	ដ							
	<b>⊢</b> – o w	4	32	ŝ	80								
	F-04	20	8	18									
	⊷ ÷ o ø	0 0	ŝ										
	0 N	4											
		101	102	103	10.1	105	106	107	108	109	C I I	111	112

TABLE 3

FREQUENCY DISTRIBUTION GPAPH FOR PAIRS OF PROGRAMS

.



TABLE 4

THE FOLLOWING PAIRS HAVE A CORRELATION OF 32: 1102,1105 1103,1106 1104,1107 TABLE 5

.

an An an an .