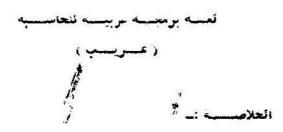
AN ARABIC COMPLIER PROGRAMMING LANGUAGE GAREB 1

Dr. Mohammaq Zeki Mohammaq * & Sanar Abqil Aziz**



نغرض بقدي نبود بحسبات بعيد كدير من يو طنين الغرب الذين لاينفنول النعه لانكسراه من صرورور البيدة بتركيب نغة برمعه عراسه للحاسبة الكول سبطه للمبتدئين وسهله التنفيد على المداك فقد تمت هذا بحارته لابناج عه مسالهه للغه بيسك) ومكتوله للغم بهرراز السلهوا العلية المراجات للخرى وقد سميت بعيبة المراجات العلية الخرى العام للمجمعة مسل وقد سميت بعيبة المراجات العلية المراجات العام العام العام المدالة المراجة عراسة للحاسبة القد العدال الحارث حاسبة المحاسبة المدالة المحاسبة المدالة المدالة عراسة المحاسبة العدالة العدالة عراسة المحاسبة العدالة المحاسبة العدالة المدالة المدال

تذالت فصيد آداب عنان صعورات لابناج هيباده اللغة ضعور حجور الدائر درمع ضد العليد آداب بالجعه أي مكانية كنانه برامج سننصه لابراند عدد حبيبا عن 30 حملة ريحتوي طفاله عالى مبال لذلك ، ربار هذه سبراي معاولة متواضعية لابناج عه عراسه للعاسلة باللي أن سطور البلد العصول على حاسلة آسرافي خامعة للاصلال للسندس البرانية ،

1- introduction :

Print of in Arabic collectors is available in many compute installations in the Arab world. However, programmers still resort to English in programming and also use than instruction media to the wide increasity to the introduction of principles of the puters to a large cross section of the Arab population in ose English is poor, it would appear that a simple programming language in A soic would be of great value.

There has been few attempts to produce compilers for arable programming languages, e.g. 4. Khawarizmi(1) Lapth(2), and others Those attempts and the present one are my at their starting and the larget of producing efficient compilers for simple arable programming languages still need extensive research

An attempt to introduce a simple anguage is presented by the authors. Two main conditions were kept in mind while starting this experiment.

First v the language should be simple and suitable for beginners. Second v a should be possible to apply it on a large number of computer insia. Mons This is why the suggested language is similar to EASIC for simplicity and coded in FORTRAN for portability.

The tackage is called (GAREB) as the reverse of tetters collected from the title. Arabic computer programming language I and its Arabic transliteration which means strange or strange or stranger and is collected from the translation of the same title.

2- Main Characteristics

The set of programs forming the FORTRAN compiler to analyse Arabic programs, should be characterised by the following

- a. The compiler should take into account the fact that arabic programs should a written from right to left.
- be permitted The blanks are a be cancelled by the compiler

^{*} Dire for of comits, femilientre, Mils il un versity

^{**} Assistant Enginees in the same is the

APPENDIX I

- (c) The Arabic characters set should take a standardised form. The set of N.C.C. - Baghdad(8) was taken as basis. It is shown in Appendix I. The keyboard layout for the data entry is shown in Appendix II.
- (d) A limited number of Arabic commands were tried because of shortage of core availability. Those are shown in Appendix III.
 - (e) Number of Arabic statements is limited mainly by the core size. For the 8K IBM 1130, it was about 20 statements. In larger installation, mor sizable programs can obviously be executed using the same compiler.

3- Bakus Normal Form (BNF) :

To define the suggested language grammatically, the Bakus Normal Form is used. The set of grammar rules of BNF has a simple structure which can be used to define the syntax of the languages. The symbol stands for " all the syntax items enclosed by brackets should be repeated". the corner brackets are used to enclose Meta-Components, and the symbol means "or " With these symbols, the syntax of the suggested language may be put as follows."

(MEAD STATEMENT) = READ (VARIABLE NAME) (VARIABLE NAME)

WHITE STATEMENT) = WRITE (VARIABLE NAME) (VARIABLE NAME)

PRITTHMETIC STATEMENT) = LET (VARIABLE NAME) = EXPRESSION)

GO TO STATEMENT) = GO TO (LABEL)

VARIABLE NAME) = SYMBOL > DIGIT) (SYMBOL) (SYMBOL)

(EXPRESSION) = (OPERAND) (OPERAND) (OPERAND)

SUBEXPRESSION) = ((EXPRESSION))

(LABEL) := (INTEGER)

(DATA STATEMENT) = DATA (CONSTANT) (CONSTANT)

OPERAND := SIGN) (DECIMAL) | SIGN) (INTEGER)

OPERATON) := SIGN) (DECIMAL) | SIGN) (INTEGER)

(PRACTION) := SIGN) (DECIMAL) | SIGN) (INTEGER)

(INTEGER) := CIGIT) (DIGIT)

(DIGIT) = 0 2 2 4 4 5 6 7 8 8

(SYMBOL) := A 5 6 C D 5 1 F G H 1 1 1 1 K L M N 1 O P M SIGN) = 1 K L M N 1 O P M SIGN) = 1 K L M N 1 O P M SIGN

M.C.C. - Baghdad ARABIC CHARACTER SET

LATIN STD	CODE	EBCDIC	ASCII V <u>ALUE</u>	ARABIC REPLACEMENT
A	11-8-7	FF	5E	۶
	0-8-5	6D	.1	ی
	9-1	79	60	1
a	12-0-1	ST	61	4
ь	12-0-	82	32	3
c	12-0-3	83	63	5
ď	12-0-4	Bla	64	<u>ج</u>
e	12-0-5	85	5.5	ح.
£	12-0-6	95	56	خ
2	12-0-7	87	67	د ڏ
g h	12-0-5	88	58	3
4	_2 - 0-J	89	59	٠
:5 :3	12-11-1	91	6A	j
k	12 - 11 - 1	92	ъв	س.
1	_2-13-	93	60	ھ
m,	_4	94	60	40
n	.2-11-1	95	52	ض
	_7-11-F	96	5F	وز
•	12-11-	97	"C.	ف
p	1,-11-	96	-1	£
Ç F	2-1	99	72	<u> 1</u>
	_i_0-0-2	A2	72	ف.
s	(*:-	EA	74	ف
		AH	75	2
u	;-	A5	to.	= 0
v	0-	AG	2.5	
W	v _1-0-	A7	79	J
×	**-^*-	AB	70	و
y	77.40-	A9	7A	<u>ء</u> عـ
2	12-1	CO	78	-
j. V	. 11-0	DO	מד	

4- Procedure of Operation :

The FORTRAN compiler performs the following operations.

- 4.1 Entering all the Arabic character set in Al-Format.

 As there are no Arabic symbols on the card punch keyboard, their latin equivalent of Appendix I was atroduced. The 47 elements comprising of alphabetic numeric and operators etc. were stored in LT.
- 4.2 Entering the Arab rommands given in Appendix III to be stored in LK.
- 4.5 Entering the Arabic program which is punched on data cards of 80 columns. It is read as Al-Format and stored in LR. The program is listed as punched on the data cards.

the array LI since the convention of formation of numbers out of digits in Arabic is the same as English. The statements at this stage can be listed in their correct form on the printer

- 4.4.4. The type of the statement using the subroutine TYPE is found by giving a code for each type according to the NTYP given in Appendix III.

 This subroutine also computes the numerical values in I-Format for the statement numbers given in A-Format using the FUNCTION GET.
- 4.5 LEXICAL ANALYSIS :(7) The most basic phase of any translation is that in which the input program is subdivided into its elementray constituents identifiers operators, symbols, numbers, ets. This phase is termed lexical analysis, and the basic program units which result from lexical analysis are termed lexical items. Typically the lexical analyzer LCLAS given in Appendix Iv is the input routine for the translator, reading successive lines of input program, breaking them down into individual lexical items, and feeding these lexical items to the later stages of the translator to be used in the

Appendix IV

Coding for Parsing

Code (LCLAS)	Object		
1	Numeric Characters		
2	Alphabetic Characters		
3	Operators		
4	Decimal point		
5	Equal sign		
6	Comma		
7	Prefix (minus or plus)		
8	Left parenthesis		
9	Right parenthesis		

higher level of analysis. Thus the texical analyzer is a translator whose input is the string of symbols representing the source program and whose output is a stream of jexical items. This output forms the input to the syntactic analyzer.

- 4.6 SYNTACTIC ANALYSIS (PARSING): The second stage in translation after the lexical analysis. Here the larger program structures are identified: statements, expressions etc. using the lexical items produced by the lexical analyzer 'LCLAS'). Syntactic analysis usually allernates with sementic analysis (4.7). First the syntactic analyzer identifies a sequence of lexical items forming a syntactic unit such as an expression or statement using subroutine LHKVR. A semantic analyzer is then called to process this unit the syntactic analyzer enters in a tack the various elements of the syntactic unit found, and these are retrieved and processed by the semantic analyzer.
- 4.7 SEMANTIC ANALYSIS: Semantic analysis is perhaps the central phase of translation of the syntactic structures recognized by the syntactic analyzer are processed and the structure of the executable object code begins to take shape.

Semantic analysis is thus the bridge between the analysis and synthesis parts of trestion. In number of other important functions also occur in this stage, including symbol table maintenance most error detection. The output from his stage is some internal form of the final executable magram.

The semantic analyzer is ordinarily split into a set of smaller semantic analyzers given in Appendix III as syntax subroutines, each of which

- 4.7.1 Handles one particular type of program construct. For example, erithmetic expressions might be handled by one analyzer. GO TO statements by another, and READ statements by another.
- 4.7.2 The appropriate semantic analyzer is called by the syntactic unalyzer whenever it has recognized a syntactic unit to be processed.
- 4.7.3 The semantic analyzer interact among themselves through information stored in various data structures, particularly in the central symbol table.
- 4.7.4 These -doroutines contain the necessary flags for

syntax errors giving a root for each type of error with resolvence to the and number where it DS LLIBO

- 4.7.5 I' the statement is correct syntactically then the speration of Tubies construction (extraction peration . Priormec The following tables are is be pristed so
 - tay Table 'vo containing see numbers in the same sequence as they come in the program
 - (b) Table NTYP containing the code numbers for the statements.
 - Table 17, containing the nature of the contehis of each statement
 - (d) Tagle MC containing the limit for each statement

By the construction of these tables the parsing phase is finished Hence the program is ready for interpretation. The abit to operation should be repeated for all the Arabic program cards

- 4.8 CODE GENERATION Trus ocoss involves formatin, the output properly ... in the information contained in the internal program representation in this stage the SYMBOL TAR E which contains information about each variable well be built The variables are defined in the MWBOL TABLE before its usage in the program.
- 4.8.1 Codes are generated by the subroutine VRTBL which tables the vriables from the read stateme nts in the program and hads them in the array VAR Simultaneous v the corresponding value of that variable is taken from the data statement using FUNCTION Co. 1 and converted to executaone numbers. This value is loaded in a table cal-I VAL Before loading any variable, it should be enecked whether was care exists in the fable or not. The table then contains the variables in all the read statements and their values

A pointer named KJ, gives the limit of the variables. contained in the Jean Statemen, from the table, Tris table is completed by adding the structures coming from the take left-hand side of the assignment statement. This is done using (STS), which adds the variables to VAR keeps g corresponding values in VAL blank this the mathematic operations are performed then their values are sucs. salet. The NTY? done is to be scanned ment statement. If any non-executable number for a are found within the statemen, sien shey are converted into executable elements using F CTION CET. This is performed using the subroutione CN VRT.

Subsouring SOLVE is Called to perform mathematical operations within assignment statemnets. The assignment operation produces a side effect-achange in the value of a variable or data structure element. The side effects produced by one issement at its the input of the next statement in sequence("

For this reason, the NTYP table is to be scanned for assignment statements in the same order as the source program The assignment statement is then moved to a new table called MATH which is used as a working table to do mathematical operations after substituting the value of the variables form the symbol table. This will save the variables in the assianment statements when it is used for more than one reration if a DO, or GOTO statements (; e) control statements are used within the program. Now, MATH consists of operands in executable form separated by operators and ready for code generation using RE-OKDERING technique which searches for the operator in this order *, +, -) When the operator is found, (t. / it become easier to find the adjecent 2 - operands (the operand, preceding and the operand following the operator i and do the mathematical operation according to the operator found(4). Then it assigns to the result a variable name within the table and reorder the table and continues in the same fashion for all the operators in the statement The final result for the assignment statement is then moved to the SYMBOL TABLE.

5. Discussions and Conclusions .

51 Limitations

Appandix V shows a solved example using the compiler described in this paper. There were few limitat ions on writing such a program such as the following

- 511 Maximum number of statements is 20 in the whole program
- 51.2 Marimum number of statements of each command $ty_1 \circ is \, 5$
- 5.1.3 Variables may consist of one or two characters .
- 514 The FORIRAN programming available on IBM 1.30 makes it difficult to convert variables from A-Format to A2, A3 ect if such a facility is a allable on any system, ther, use of polish form