

Εργασία Μεταγλωττιστών Τμήμα Α2

Ομάδα 15

29 Μαρτίου 2021



Ομάδα 15
Αναλυτικά τα μέλη:

Διονύσης Νικολόπουλος *AM* : 18390126
Θανάσης Αναγνωστόπουλος *AM* : 18390043
Αριστείδης Αναγνωστόπουλος *AM* : 16124
Σπυρίδων Φλώρος *AM* : 141084

Αναλυτικά οι ρόλοι:

Γενικός Συντονιστής: Διονύσης Νικολόπουλος
Υπεύθυνος Τμήματος Εργασίας Α2: Διονύσης Νικολόπουλος

Η εργασία αυτή πραγματοποιήθηκε με χρήση L^AT_EX

Περιεχόμενα

1	Εισαγωγή	2
2	Περιγραφή σύνταξης της Uni-C σε BNF	4
3	Διαγράμματα Ενιαίου Αυτόματου	5
3.1	Γενικευμένο Διάγραμμα Ενιαίου Αυτόματου	5
3.2	Λεπτομερές Διάγραμμα Ενιαίου Αυτόματου	6
4	Εξήγηση Κώδικα σε FSM	7
5	Ολοκληρωμένος Κώδικας σε FSM	13

1 Εισαγωγή

Τι πραγματεύεται η εργασία

Στην εργαστηριακή εργασία αυτή ασχοληθήκαμε με την δημιουργία ενός ντετερμινιστικού αυτόματου πεπερασμένων καταστάσεων.

Ο σκοπός του αυτομάτου

Το αυτόματο αυτό έχει ως σκοπό να λειτουργεί σαν Λεκτικός Αναλυτής της γλώσσας Uni-C, ενός παρακλαδιού της γλώσσας C, απο το Πανεπιστήμιο Δυτικής Αττικής.

Βήματα Υλοποίησης

Πριν πραγματοποιηθεί φυσικά η κωδικοποίηση για την προσέγγιση του προβλήματος αυτού, έγινε η περιγραφή της γλώσσας Uni-C σε μορφή BNF.

Η κωδικοποίηση του έχει γίνει καταρχήν με την βοήθεια του προγράμματος *FSM*, που επίσης είναι προγραμματισμένο πάνω στην γλώσσα C.

Ο σχεδιασμός των διαγραμμάτων έχει γίνει με το πρόγραμμα *JFLAP*.

Εξήγηση για το τι ακολουθεί

Σε αυτό το PDF αρχικά θα παρατεθεί η περιγραφή της Uni-C σε BNF.

Έπειτα, θα παρατεθούν επίσης και τα διαγράμματα του εναίου αυτόματου που βοηθούν στην κωδικοποίηση του αλλά και στην σύλληψη νοητά του τι περιμένουμε να κάνει το αυτόμάτο μας, πως θα το υλοποιήσουμε τεχνικά.

Στην συνέχεια, θα ακολουθήσει ο κώδικας του αυτόματου τμηματικά με εξηγήσεις της λειτουργίας του.

Τελος, εφόσον ο κώδικας παραδωθεί και ολοκληρωμένα, θα παρατεθούν ορισμένοι ελέγχοι πραγματικής λειτουργίας του αυτόματου σε εικόνες.

2 Περιγραφή σύνταξης της Uni-C σε BNF

Πριν την υλοποίηση του αυτομάτου σε εκτελέσιμο κώδικα, όπως αναφέραμε παραπάνω, έγινε η ανάλυση της σύνταξης της γλώσσας Uni-C σε μορφή BNF (Backus-Naur Form ή Μορφή Μπάκους-Ναούρ).

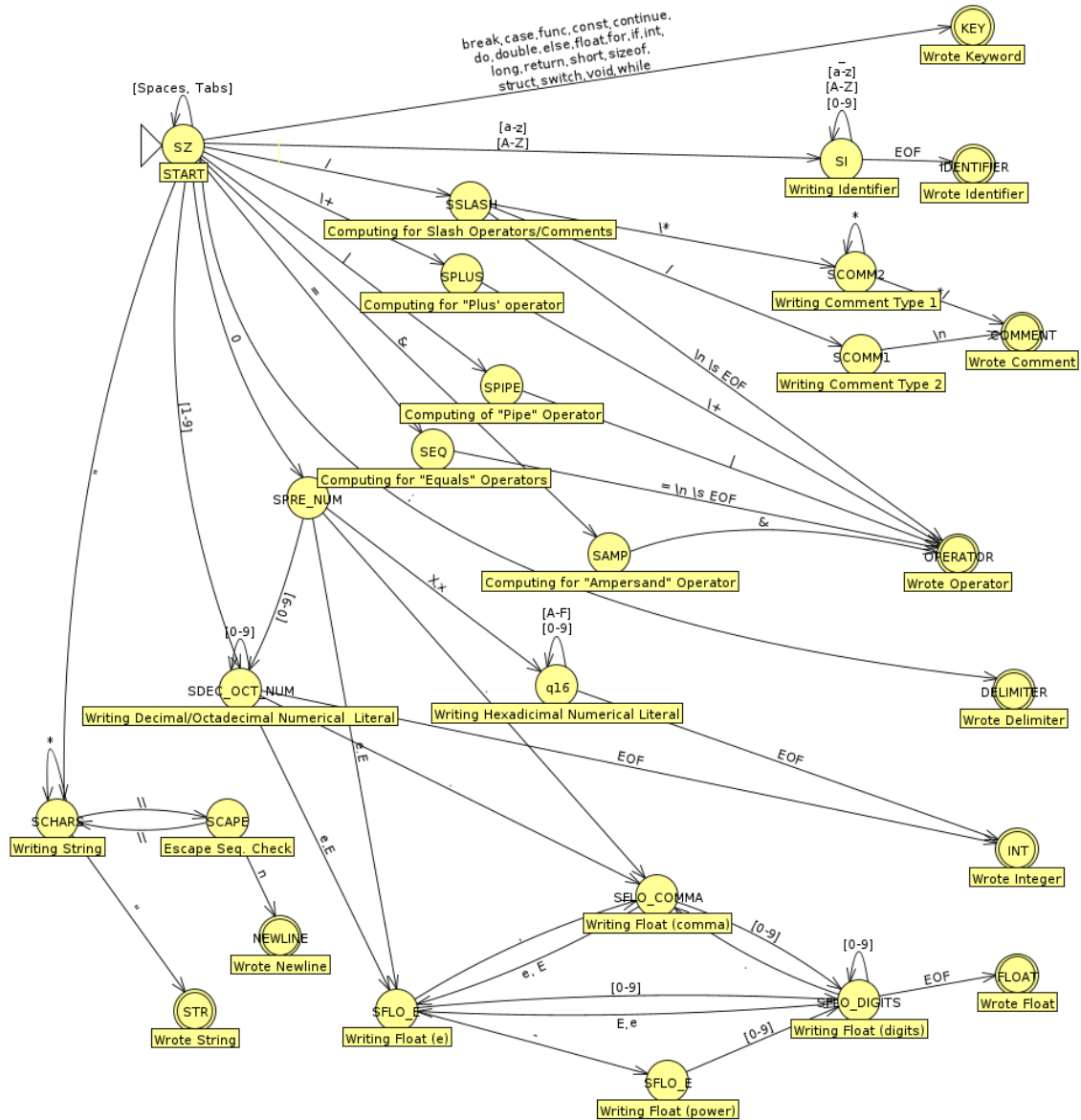
Οπότε, με την βοήθεια του περιβάλλοντος λεκτικού επεξεργαστή Emacs, γράψαμε τον εξής συμβολισμό.

```
1 symbol      ::= "|" | " " | "!" | "#" | "$" | "%" | "&" | "(" | ")"
  " | "*" | "+" | "," | "-" | "." | "/" | ":" | ";" | ">" | "=" |
  "<" | "?" | "@" | "[" | "\" | "]" | "^" | "_" | "`" | "{" | "}" |
  "~"
2 letter      ::= "A" | "B" | "C" | "D" | "E" | "F" | "G" | "H" | "I"
  " | "J" | "K" | "L" | "M" | "N" | "O" | "P" | "Q" | "R" | "S" | "
  T" | "U" | "V" | "W" | "X" | "Y" | "Z" | "a" | "b" | "c" | "d" |
  "e" | "f" | "g" | "h" | "i" | "j" | "k" | "l" | "m" | "n" | "o" |
  "p" | "q" | "r" | "s" | "t" | "u" | "v" | "w" | "x" | "y" | "z"
3 digit       ::= "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" |
  "8" | "9"
4 character   ::= letter | symbol | digit
5 character1   ::= character | "'"
6 character2   ::= character | ','
7 text1       ::= " | character1 text1
8 text2       ::= ' | character2 text2
9 literal     ::= ''' text1 ''' | '"' text2 '"'
10 identifier  ::= character text1
11 operator    ::= "+" | "-" | "*" | "/" | "%" | "=" | "+=" | "-=" |
  "*=" | "" | "" | "--" | "=" | "=" | "&" | "==" | "!=" | "++" |
  "&&" | "||"
12 number1    ::= digit number1 | "'"
13 number2    ::= digit number2 | ','
14 number     ::= number1 | number2
15 integer    ::= number | "0X" number | "0x" number | "0"
16 float      ::= number "." number | number "e" number | number "e"
  number "-" number
17
```


3.2 Λεπτομερές Διάγραμμα Ενιαίου Αυτόματου

Εδώ είναι το λεπτομερές ενιαίο αυτόματο που σχεδιάσαμε:

Σχήμα 2: Λεπτομερές Αυτόματο



4 Εξήγηση Κώδικα σε FSM

Το πρώτο κομμάτι κώδικα σε FSM είναι, ασφαλώς, η αρχική μας κατάσταση:

```
1 START=SZ
2 SZ: % -> OPERATOR
3     ! < > \- = \* -> SEQ
4     \+ -> SPLUS
5     & -> SAMP
6     ; -> DELIMITER
7     | -> SPIPE
8     a-z A-Z _ -> SI
9     \n \s -> SZ
10    / -> SSLASH
11    " -> SCHARS
12    0 -> SPRE_NUM
13    1-9 -> SDEC_OCT_NUM
14    * -> BAD
15
```

Εδώ, βλέπουμε ότι ο αναλυτής μας περιμένει τον πρώτο, ουσιαστικά χαρακτήρα από τον χρήστη, ο οποίος θα έχει μεγάλη σημασία στο να καθορίσουμε "τι πάει να γράψει ο χρήστης".

Βλέπουμε πολλές καταστάσεις, των οποίων η φιλοσοφία εξηγείται ως εξής.

- Οι καταστάσεις που το όνομα τους αρχίζει με "S" είναι για τον αναλυτή μας οι "μεταβατικές" μας καταστάσεις, οι καταστάσεις στις οποίες δεν είναι απόλυτα σίγουρο το τι θέλει να γράψει ο χρήστης, αλλά είναι σίγουροι, με βάση κανόνες και ντετερμινιστικά οι "δρόμοι" που μπορεί να ακολουθήσει ο χρήστης.
- Οι υπόλοιπες καταστάσεις, καταστάσεις όπως η DELIMITER που βλέπουμε εδώ, εκφράζουν επιστροφή του αναγνωριστικού (και άρα επιτυχής αναγνώριση) της λεκτικής μονάδας (πχ εδώ το token (λεκτική μονάδα) είναι DELIMETER) καθώς ο αναλυτής μας είναι σίγουρος, με βάση κανόνες που ακολουθεί ντετερμινιστικά, ότι η λέξη η οποία ο χρήστης εισήγαγε είναι ενός συγκεκριμένου τύπου.

Συνεχίζοντας, βλέπουμε 2 καταστάσεις που αφορούν την αναγνώριση συμβολοσειράς, στις οποίες είναι δυνατόν να μεταβεί το αυτόματο, αν ο χρήστης εισάγει σε ορισμένες καταστάσεις τον χαρακτήρα `"`.

Αυτό ισχύει για την αρχική κατάσταση, όπως είδαμε παραπάνω, αλλά και για τις καταστάσεις IDENTIFIER, INT, FLOAT, αφού είναι πιθανό ο χρήστης να γράφει μια συμβολοσειρά έπειτα από τους τύπους λέξεων αυτούς.

Η κατάσταση SCHARS συγκεκριμένα είναι αυτή στην οποία δίνονται οι χαρακτήρες της συμβολοσειράς, και αφήνεται ένα ενδεχόμενο είτε να τερματισθεί η συμβολοσειρά με το να εισαχθεί ξανά ο χαρακτήρας `"` (και να μεταβεί στην κατάσταση STR), είτε να εισαχθεί ο χαρακτήρας `\`.

Στην τελευταία περίπτωση, ο αναλυτής έρχεται στην κατάσταση SCAPE, την οποία περιμένει να εισαχθεί είτε ο χαρακτήρας `n`, οπότε να επιστραφεί η κατάσταση NEWLINE, είτε να εισαχθεί ο χαρακτήρας `\` ξανά, οπότε να επιστρέψει στην κατάσταση SCHARS, είτε να κάνει λάθος και να εισαχθεί οποιοσδήποτε άλλος χαρακτήρας, οπότε να πάει στην κατάσταση BAD.

Στην κατάσταση NEWLINE, οποιοσδήποτε χαρακτήρας μας πάει πίσω στην κατάσταση SCHARS, ενώ ο χαρακτήρας `"` στην κατάσταση STR

```
1 SCHARS: * -> SCHARS
2         " -> STR
3         \\ -> SCAPE
4
5 SCAPE:  n -> NEWLINE
6         \\ -> SCHARS
7         * -> BAD
8
9 NEWLINE: " -> STRING
10        * -> SCHARS
11
12 STR:    \s \n -> SZ
13
```

Με παρόμοια λογική βλέπουμε διάφορους χαρακτήρες που είναι *δυνητικά* τελεστές (ο πρώτος είναι *δυνητικά* και σχόλιο).

```
1 SSLASH: / -> SCOMM1
2         = \n \s -> OPERATOR
3         \* -> SCOMM2
4         * -> BAD
5
6 SEQ:    = \n \s -> OPERATOR
7         * -> BAD
8
9 SPLUS:  = \+ -> OPERATOR
10        * -> BAD
11
12 SAMP:   & \s -> OPERATOR
13        * -> BAD
14
15 SPIPE:  | -> OPERATOR
16        * -> BAD
17
```

Επίσης, εδώ βλέπουμε, με παρόμοιο συλλογισμό με τις καταστάσεις για τις συμβολοσειρές, και την επεξεργασία για τα σχόλια (με τελική κατάσταση την κατάσταση COMMENT)

```
1 SCOMM1:  * -> SCOMM1
2          \n -> SZ
3
4 SCOMM2:  \* -> SCOMM21
5          * -> SCOMM2
6
7 SCOMM21: / -> COMMENT
8          * -> SCOMM2
9
10 COMMENT: \n \s -> SZ
11
```

Παρατηρούμε και την κατάσταση SI, η οποία δέχεται τους επιτρεπτούς χαρακτήρες ως κομμάτια της πιθανής λέξης ενός identifier, που με τον χαρακτήρα `;`, με κενό, ή με τέλος αρχείου γίνεται η μεταβολή του αναλυτή στην κατάσταση IDENTIFIER.

```
1 SI:      a-z A-Z 0-9 -> SI
2          ; -> IDENTIFIER
3          \s -> IDENTIFIER
4          EOF -> IDENTIFIER
```

```
5      * -> BAD
```

```
6
```

Εδώ βλέπουμε τις 'τελικές' καταστάσεις στις οποίες φτάνει ο αναλυτής όταν κρίνει μια λέξη να πληρεί τις προϋποθέσεις για να χαρακτηριστεί ως μια συγκεκριμένη λεκτική μονάδα.

Παρατηρούμε ότι, με βάση την γλώσσα Uni-C, οι λεκτικές μονάδες είναι δυνατό να ακολουθηθούν από συγκεκριμένες λεκτικές μονάδες (παράδειγμα οι τελεστές μονάχα από ένα νούμερο ή identifier, για αυτό οι μοναδικές επιτρεπτές καταστάσεις που μεταβαίνει ο αναλυτής μετά από την κατάσταση OPERATOR είναι η ενδιαμέση κατάσταση είτε αριθμού είτε identifier)

```
1  OPERATOR:   a-z A-Z _ -> SI
2              0 -> SPRE_NUM
3              1-9 -> SDEC_OCT_NUM
4              " -> SCHARS
5              * -> BAD
6
7  DELIMITER:  * -> SZ
8  INT:        \s \n -> SZ
9              A-Z a-z _ -> SI
10             % -> OPERATOR
11             ! < > \- = \* -> SEQ
12             \+ -> SPLUS
13             & -> SAMP
14             ; -> DELIMITER
15             | -> SPIPE
16             / -> SSLASH
17             0 -> SPRE_NUM
18             1-9 -> SDEC_OCT_NUM
19             * -> BAD
20             " -> SCHARS
21
22  FLOAT:      \s \n -> SZ
23             A-Z a-z _ -> SI
24             % -> OPERATOR
25             ! < > \- = \* -> SEQ
26             \+ -> SPLUS
27             & -> SAMP
28             ; -> DELIMITER
29             | -> SPIPE
30             / -> SSLASH
31             0 -> SPRE_NUM
32             1-9 -> SDEC_OCT_NUM
```

```
33          * -> BAD
34
35 IDENTIFIER: \s \n -> SZ
36          A-Z a-z _ -> SI
37          % -> OPERATOR
38          ! < > \- = \* -> SEQ
39          \+ -> SPLUS
40          & -> SAMP
41          ; -> DELIMITER
42          | -> SPIPE
43          / -> SSLASH
44          0 -> SPRE_NUM
45          1-9 -> SDEC_OCT_NUM
46          * -> BAD
47          " -> SCHARS
48
49
50
```

Τέλος, έχουμε τις ενδιαμέσες καταστάσεις με τις οποίες ο αναλυτής κρίνει αν ο αριθμός που εισάγεται είναι πραγματικός ή ακέραιος.

Με βάση λοιπόν την σύνταξη της Uni-C, διαχωρίζονται αυτές οι δύο κατηγορίες.

(πχ. ακέραιοι: 0, 954, 4235 πχ. πραγματικοί: 3.14 2.43E31, 35.3e10, 34e12, 0e0, 0.2222)

```
1 SPRE_NUM:    x X -> SHEX
2              0-9 -> SDEC_OCT_NUM
3              \. -> SFLO1
4              e E -> SFLO2
5              * -> BAD
6
7 SDEC_OCT_NUM: 0-9      -> SDEC_OCT_NUM
8              \.      -> SFLO_COMMA
9              e E      -> SFLO_E
10             \s \n ; -> FLOAT
11             * -> BAD
12
13 SFLO_COMMA:   \s \n ; -> FLOAT
14             e E      -> SFLO_E
15             0-9      -> SFLO_DIGITS
16             * -> BAD
17
18 SFLO_DIGITS:  0-9      -> SFLO_DIGITS
19             \n \s ; -> FLOAT
20             e E      -> SFLO_E
21             \.      -> SFLO_COMMA
22             * -> BAD
23
24 SFLO_E:       \.      -> SFLO_COMMA
25             \-      -> SFLO_POWER
26             0-9      -> SFLO_DIGITS
27             * -> BAD
28
29 SFLO_POWER:   0-9      -> SFLO_DIGITS
30             * -> BAD
31
32 SHEX:         A-F      -> SHEX
33             0-9      -> SHEX
34             \s \n ; -> INT
35             * -> BAD
36
```

5 Ολοκληρωμένος Κώδικας σε FSM

Ο κώδικας του αυτοματού μας μπορεί να βρεθεί ως αρχείο και στον φάκελο `Source_Code` που βρήκεται μέσα στον κύριο φάκελο του παραδοτέου.

Ακολουθεί ο ολοκληρωμένος κώδικας σε FSM του αυτόματου μας:

```
1 START=SZ
2 SZ: % -> OPERATOR
3     ! < > \- = \* -> SEQ
4     \+ -> SPLUS
5     & -> SAMP
6     ; -> DELIMITER
7     | -> SPIPE
8     a-z A-Z _ -> SI
9     \n \s -> SZ
10    / -> SSLASH
11    " -> SCHARS
12    0 -> SPRE_NUM
13    1-9 -> SDEC_OCT_NUM
14    * -> BAD
15
16
17 SCHARS: * -> SCHARS
18         " -> STR
19         \\ -> SCAPE
20
21 SCAPE:  n -> NEWLINE
22         \\ -> SCHARS
23         * -> BAD
24
25 NEWLINE: " -> STRING
26          * -> SCHARS
27
28 SSLASH: / -> SCOMM1
29         = \n \s -> OPERATOR
30         \* -> SCOMM2
31
32 SEQ:    = \n \s -> OPERATOR
33         * -> BAD
34
35 SPLUS:  = \+ -> OPERATOR
36         * -> BAD
37
38 SAMP:    & \s -> OPERATOR
```

```

39      * -> BAD
40
41 SPIPE:  |      -> OPERATOR
42      * -> BAD
43
44 SI:      a-z A-Z 0-9 -> SI
45          ;          -> IDENTIFIER
46          \s         -> IDENTIFIER
47          EOF        -> IDENTIFIER
48      * -> BAD
49
50 SCOMM1:  * -> SCOMM1
51          \n -> SZ
52
53 SCOMM2:  \* -> SCOMM21
54          * -> SCOMM2
55
56 SCOMM21: / -> COMMENT
57          * -> SCOMM2
58
59 COMMENT: \n \s -> SZ
60
61 OPERATOR: a-z A-Z _ -> SI
62           0 -> SPRE_NUM
63           1-9 -> SDEC_OCT_NUM
64           " -> SCHARS
65           * -> BAD
66
67 DELIMITER: * -> SZ
68
69 SPRE_NUM:  x X -> SHEX
70           0-9 -> SDEC_OCT_NUM
71           \. -> SFLO1
72           e E -> SFLO2
73           * -> BAD
74
75 SDEC_OCT_NUM: 0-9      -> SDEC_OCT_NUM
76              \.      -> SFLO_COMMA
77              e E      -> SFLO_E
78              \s \n ; -> FLOAT
79              * -> BAD
80
81 SFLO_COMMA:  \s \n ; -> FLOAT
82              e E      -> SFLO_E

```



```

83          0-9      -> SFLO_DIGITS
84          * -> BAD
85
86 SFLO_DIGITS:  0-9      -> SFLO_DIGITS
87              \n \s ; -> FLOAT
88              e E      -> SFLO_E
89              \.      -> SFLO_COMMA
90              * -> BAD
91
92 SFLO_E:       \.      -> SFLO_COMMA
93              \-      -> SFLO_POWER
94              0-9      -> SFLO_DIGITS
95              * -> BAD
96
97 SFLO_POWER:   0-9      -> SFLO_DIGITS
98              * -> BAD
99
100 SHEX:        A-F      -> SHEX
101              0-9      -> SHEX
102              \s \n ; -> INT
103              * -> BAD
104
105 INT:          \s \n -> SZ
106              A-Z a-z _ -> SI
107              % -> OPERATOR
108              ! < > \- = \* -> SEQ
109              \+ -> SPLUS
110              & -> SAMP
111              ; -> DELIMITER
112              | -> SPIPE
113              / -> SSLASH
114              0 -> SPRE_NUM
115              1-9 -> SDEC_OCT_NUM
116              * -> BAD
117              " -> SCHARS
118
119 FLOAT:        \s \n -> SZ
120              A-Z a-z _ -> SI
121              % -> OPERATOR
122              ! < > \- = \* -> SEQ
123              \+ -> SPLUS
124              & -> SAMP
125              ; -> DELIMITER
126              | -> SPIPE

```

```
27         / -> SSLASH
28         0 -> SPRE_NUM
29         1-9 -> SDEC_OCT_NUM
30         * -> BAD
31
32 IDENTIFIER: \s \n -> SZ
33         A-Z a-z _ -> SI
34         % -> OPERATOR
35         ! < > \- = \* -> SEQ
36         \+ -> SPLUS
37         & -> SAMP
38         ; -> DELIMITER
39         | -> SPIPE
40         / -> SSLASH
41         0 -> SPRE_NUM
42         1-9 -> SDEC_OCT_NUM
43         * -> BAD
44         " -> SCHARS
45
46 STR:      \s \n -> SZ
47
48 BAD(OK): \n -> SZ
49 GOOD(OK): \n -> GOOD
50
51
```