${\bf Trabajo\ Bioinform\'atica}^*$

Ignacio Amat Hernández

Ángela LLopis Hernández

Estrella Ballester Hoyo

1 de mayo de $2020\,$

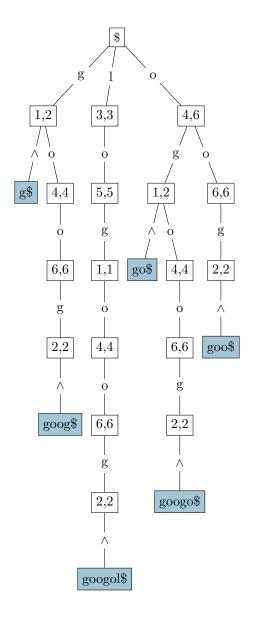


Figura 1: Árbol de prefijos con intervalos SA de la palabra "googol\$".

^{*}Grado en Ingeniería Biomédica, Escuela Técnica Superior de Ingenieros Industriales, Valencia, España.

1. Materiales

Los lenguajes que hemos comparado son \mathbb{R} , python y \mathbb{C} . La implementación en \mathbb{R} es la misma a la usada en las prácticas de la asignatura. Con la implementación proporcionada hemos dividido dos versiones, una que dibuja la traza y otra que calcula el resultado de forma limpia. Todas las versiones del código escrito está disponible en el anejo.

El tiempo y la memoria la hemos medido de dos maneras. Primero hemos evaluado a la función únicamente usando herramientas del lenguaje. En python hemos usado los paquetes resource y time para medir el tiempo y la memoria. En hemos usado la utilidad system. time para el tiempo y la libreria profmem para perfilar el uso de memoria. En hemos medido el tiempo usando la librería time. h y la memoria usando la herramienta Valgrind. Después hemos medido el uso de recursos desde la línea de comandos, llamando a los programas como scripts usando funcionalidades propias de cada lenguaje. Con ello podemos evaluar el sobrecoste de cada una de las tres plataformas usadas. Para ello hemos empleado la herramienta bench y el comando time.

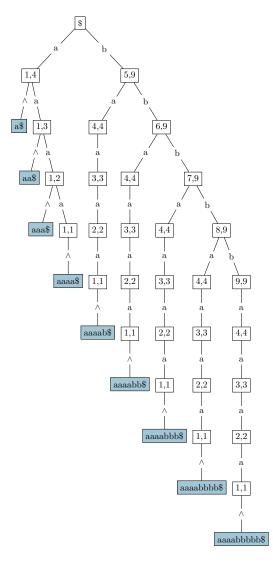


Figura 2: Árbol de prefijos con intervalos SA de la palabra "aaaabbbbb\$".

2. Resultados

Hemos implementado las trazas de manera que son perfectamente idénticas en python y en . A continuación mostramos algunos ejemplos comparativos para verificar que todos los programas se comportan de forma igual. Posteriormente comparamos el uso de tiempo y memoria.

1		Mutation	i	\mathbf{z}	k	1
2	Deletion	[1]	1 -	-1	0	6
3	Insertion	[g]	2 -	-1	1	2
4	Substitution	[g] -> [l]	1 -	-1	1	2
5	Insertion	[1]	2 -	-1	3	3
6	Match	[1]	1	0	3	3
7	Deletion	[o]	0 -	-1	3	3
8	Insertion	[o]	1 -	-1	5	5
9	Match	[o]	0	0	5	5
10	Deletion	[g]	-1	-1	5	5
11	Insertion	[g]	0 -	-1	1	1
12	Match	[g]	-1	0	1	1
13		{1,1}				
14	Insertion	[o]	2 -	-1	4	6
15	Substitution	[o] -> [1]	1 -	-1	4	6

Cuadro 1: Traza de INEXRECUR con X = "googol\$", W = "gol", z = 0 en 😉 y 🤚 python"

```
INEXRECUR
                     by XAVI GABRI AITANA ALFREDO
1
2
3
                         2
       Ι
               g ]
                             -1
4
                            -1
              [ 1 ]
5
       Ι
                         2
                                  3
                             -1
6
       Μ
                                  3
                                      3
7
           D
                  [ o ]
8
                  [ o ]
                                           5
9
                  [ o ]
                                  0
                                      5
10
                      [ g ]
               Ι
                                  0
11
                       [ g ]
12
               Μ
                       [ g ]
                                  -1
                                       0
  2
                                       6
14
                                     6
15
                            -1
```

Cuadro 2: Traza de INEXRECUR con X = "googol\$", W = "gol", z = 0 en

1		Mutation	i	z	k	1
2	Deletion	[g]	2	-1	0	6
3	Insertion	[g]	3	-1	1	2
4	Match	[g]	2	0	1	2
5	Deletion	[o]	1	-1	1	2
6	Insertion	[o]	2	-1	4	4
7	Match	[o]	1	0	4	4
8	Deletion	[o]	0	-1	4	4
9	Insertion	[o]	1	-1	6	6
10	Match	[o]	0	0	6	6
11	Insertion	[1]	3	-1	3	3
12	Substitution	[l] -> [g]	2	-1	3	3
13	Insertion	[o]	3	-1	4	6
14	Substitution	[o] -> [g]	2	-1	4	6

Cuadro 3: Traza de Inex Recur con X = "googol\$", W = "goog
", z = 0 en \bigcirc y \bigcirc python

1		INE	XREC	UR	_	by	XAV	T GABF	RI AITA	ANA AI	LFREDO)	
2	_	D		[g]		2	-1	0	6			
3	_	I		[g]		3	-1	1	2			
4	_	\mathbf{M}		[g]		2	0	1	2			
5	_	_	D		[o]		1	-1	1	2		
6	_	_	I		[o]		2	-1	4	4		
7	_	_	\mathbf{M}		[o]		1	0	4	4		
8	_	_	_	D		[o]		0	-1	4	4	
9	_	_	_	I		[o]		1	-1	6	6	
10	_	_	_	\mathbf{M}		[o]		0	0	6	6	
11	_	I		[1]		3	-1	3	3			
12	_	\mathbf{S}	[1	->	g]	2	2	-1	3	3			
13	_	I		o]]		3	-1	4	6			
14	_	\mathbf{S}	[o	->	g]	2	2	-1	4	6			

Cuadro 4: Traza de Inex
Recur con X = "googol\$", W = "goog", z = 0 en ${\bf R}$

1		Mutation	i	\mathbf{z}	k	1
2	Deletion	[1]	2	0	1	6
3	Deletion	[o]	1	-1	1	6
4	Insertion	[g]	2	-1	1	2
5	Substitution	[g] -> [o]	1	-1	1	2
6	Insertion	[o]	2	-1	4	6
7	Match	[o]	1	0	4	6
8	Deletion	[o]	0	-1	4	6
9	Insertion	[g]	1	-1	1	2
10	Substitution	[g] -> [o]	0	-1	1	2
11	Insertion	[o]	1	-1	6	6
12	Match	[o]	0	0	6	6
13	Deletion	[g]	-1	-1	6	6
14	Insertion	[g]	0	-1	2	2
15	Match	[g]	-1	0	2	2
16 -		—— {2,2} —				
17	Insertion	[g]	3	0	1	2
18	Substitution	[g] -> [l]	2	0	1	2
19	Deletion	[o]	1	-1	1	2
20	Insertion	[o]	2	-1	4	4
21	Match	[o]	1	0	4	4
22	Deletion	[o]	0	-1	4	4
23	Insertion	[o]	1	-1	6	6
24	Match	[o]	0	0	6	6
25	Deletion	[g]	-1	-1	6	6
26	Insertion	[g]	0	-1	2	2
27	Match	[g]	-1	0	2	2
28 -		—— {2,2} —				
29	Insertion	[o]	3	0	4	6
30	Substitution	[o] -> [l]	2	0	4	6
31	Deletion	[o]	1	-1	4	6
32	Insertion	[g]	2	-1	1	2
33	Substitution	[g] -> [o]	1	-1	1	2
34	Insertion	[o]	2	-1	6	6
35	Match	[o]	1	0	6	6
36	Deletion	[o]	0	-1	6	6
37	Insertion	[g]	1	-1	2	2
38	Substitution	[g] -> [o]	0	-1	2	2

Cuadro 5: Traza de INEXRECUR con X = "googol\$", W = "gool", z = 1 en \bigcirc y \bigcirc python"

```
INEXRECUR – by XAVI GABRI AITANA ALFREDO
1
        [ \quad l \quad ] \qquad \qquad 2
2
                      0
                         1
                             6
     \mathbf{D}
3
        D
            [ o ]
                      1
                          -1
                             1
        Ι
             [ g ]
                      2
4
                          -1
        S
           [ g -> o ] 1
                        -1
5
                             1
                          -1
6
        Ι
            [ o ]
                       2
                              4
7
        M
             [ o ] 1
                          0
                             4
            D [ o ]
8
                          0
                                     6
                             -1
9
            Ι
               [ g ]
                         1
                             -1
            \mathbf{S}
10
               [ g -> o ]
                        0
                            -1
                                1
11
            Ι
               [ o ]
                         1
                             -1
                                 6
                                     6
               [ o ] 0
12
            Μ
                             0
                                 6
                                 -1
13
               D [ g ]
                             -1
14
               Ι
                    [ g ]
                             0
                                 -1
                                     2
15
               M
                  [ g ]
                            -1
                                 0
  16
              3 	 0 	 1
17
           [ g ]
         [ g -> l ] 2 0 1
18
          [ o ]
19
                      1
                          -1
             [ o ]
                      2
                          -1
20
        Ι
             [ o ]
21
                          0
                             4
        Μ
              [ o ]
22
            D
                          0
                             -1
                [ o ]
23
            Ι
                             -1
               [ o ]
24
            Μ
                          0
                             0
                                 6
                 [ g ]
               D
25
                             -1
                                 -1
                    [ g ]
26
               Ι
                             0
                                 -1
                                     2
                  [ g ]
                                     2
                             -1
                                 0
27
               Μ
  28
        [ o ] 3 0 4
                            6
29
     Ι
        [ o -> 1 ] 2 0 4 6
30
        D [ o ]
                      1
                         -1
31
32
        Ι
           [ g ]
                      2
                         -1
           [ g -> o ]
33
        S
                    1
                        -1
                            1
            [ o ]
34
        Ι
                      2
                             6
                          -1
35
        M
             [ o ]
                      1
                          0
                             6
            D [ o ]
36
                         0
                             -1
37
            I
                [ g ]
                         1
                             -1
              [g \rightarrow o]
                            -1
38
```

Cuadro 6: Traza de INEXRECUR con X = "googol\$", W = "gool", z = 1 en 😱

Se puede comprobar que las tres trazas contienen la misma información. Procedemos ahora a la comparación de los tiempos medidos, primero las mediciones tomadas desde el propio programa.

```
inexrecur_time.c
1
2
    12.34\mu s from 10000 iterations.
3
    inexrecur_time.py
4
5
    real
              sys
                      user
6
    268.89\mu s 0.37 \mu s 268.16 \mu s
7
8
    inexrecur_time.R
9
    user
               system
                           elapsed
   5422\mu s
10
               18\mu s
                           5443 \mu s
```

Cuadro 7: Tiempos de "CPU".

En cada caso tomamos múltiples lecturas y hacemos la media. Como es de esperar la solución compilada de es ~20 veces más rápida al script de python y parece ser ~450 veces más lento que y ~20 veces más lento que python. Para comparar el efecto del interpretador medimos ahora el tiempo ejecutando desde la línea de comandos. LLamamos a los scripts usando #!/bin/env Rscript y #!/bin/env Python.

```
1
    bench ./inexrecur_clean ./inexrecur_clean.py ./inexrecur_clean.R
    benchmarking bench/./inexrecur_clean
                               4.524~\mathrm{ms}
                                              (4.496 \text{ ms} \dots 4.551 \text{ ms})
 3
    time
                                0.999 R^2
                                               (0.999 R^2 ... 1.000 R^2)
 4
                                4.522 \text{ ms}
                                               (4.499 ms .. 4.559 ms)
 5
    mean
 6
    std dev
                                                (58.07 \ \mu s \ .. \ 133.1 \ \mu s)
                                89.83 \ \mu s
 7
 8
    benchmarking bench/./inexrecur_clean.py
                                39.26 \, \mathrm{ms}
                                              (39.12 ms .. 39.40 ms)
9
    time
10
                                1.000 R^2
                                              (1.000 R^2 \dots 1.000 R^2)
                                39.30 \text{ ms}
                                              (39.18 ms .. 39.45 ms)
11
    mean
12
    std dev
                                266.5 \ \mu s
                                              (177.4 \ \mu s \ \dots \ 388.7 \ \mu s)
13
14
    benchmarking bench/./inexrecur_clean.R
                                              (273.6 ms .. 285.0 ms)
15
    time
                                278.5 \text{ ms}
                                              (1.000 R^2 \dots 1.000 R^2)
                                1.000 R^2
16
                                280.6 \, \mathrm{ms}
                                              (279.1 ms .. 281.9 ms)
17
    mean
    std dev
                                1.714 \text{ ms}
                                              (1.027 \text{ ms} \dots 2.517 \text{ ms})
18
    variance introduced by outliers: 16% (moderately inflated)
19
```

Cuadro 8: Tiempo de "pared" según la utilidad bench.

En este caso \bigcirc es \sim 10 veces más rápido a \bigcirc python $^{\circ}$ y \sim 70 veces más rápido que \bigcirc , que continua siendo \sim 7 veces más lento que \bigcirc python $^{\circ}$.

```
inexrecur_clean.c
==81469== in use at exit: 18,588 bytes in 164 blocks
==81469== total heap usage: 191 allocs, 27 frees, 24,894 bytes allocated
inexrecur_mem.py
6,631,424 bytes
inexrecur_mem.R
9,4,932,584 bytes
```

Cuadro 9: Memoria usada medido desde el script.

La memoria usada es comparable en R y python cuando la medimos sin tener en cuenta el entorno de ejecución, utiliza entorno a 200 veces menos memoria.

```
time ./inexrecur_clean
 2
   (5,5)
 3
   (1,1)
    ./inexrecur_clean
   0.00s user 0.00s system 44% cpu 0.004 total
 6
   max memory:
                                     676 \text{ kB}
 7
    time ./inexrecur_clean.py
    (1, 1)
9
10
    (5, 5)
    ./inexrecur_clean.py
11
    0.03\,\mathrm{s} user 0.01\,\mathrm{s} system 82\,\% cpu 0.039 total
12
                                     6272 \text{ kB}
13
    max memory:
14
    time ./inexrecur_clean.R
15
16
    [[1]]
17
    [1] 5 5
18
19
    [[2]]
    [1] 1 1
20
21
    ./inexrecur_clean.R
22
23
    0.23\,\mathrm{s} user 0.04\,\mathrm{s} system 97\,\% cpu 0.277 total
    max memory:
                                     67476 kB
```

Cuadro 10: Memoria usada medido desde fuera del script.

Cuando consideramos el entorno de ejecución interpretado al completo \bigcirc se mantiene a la cabeza con la sorpresa de que \bigcirc emplea ~ 10 veces más memoria que \bigcirc python en este ejemplo en particular.

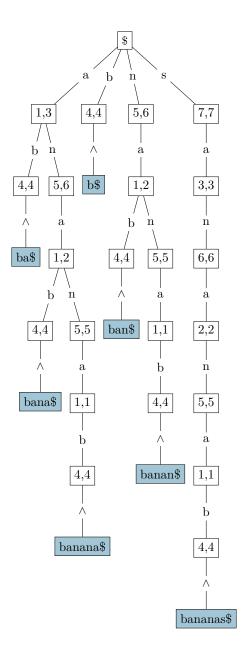


Figura 3: Árbol de prefijos con intervalos SA de la palabra "bananas\$".

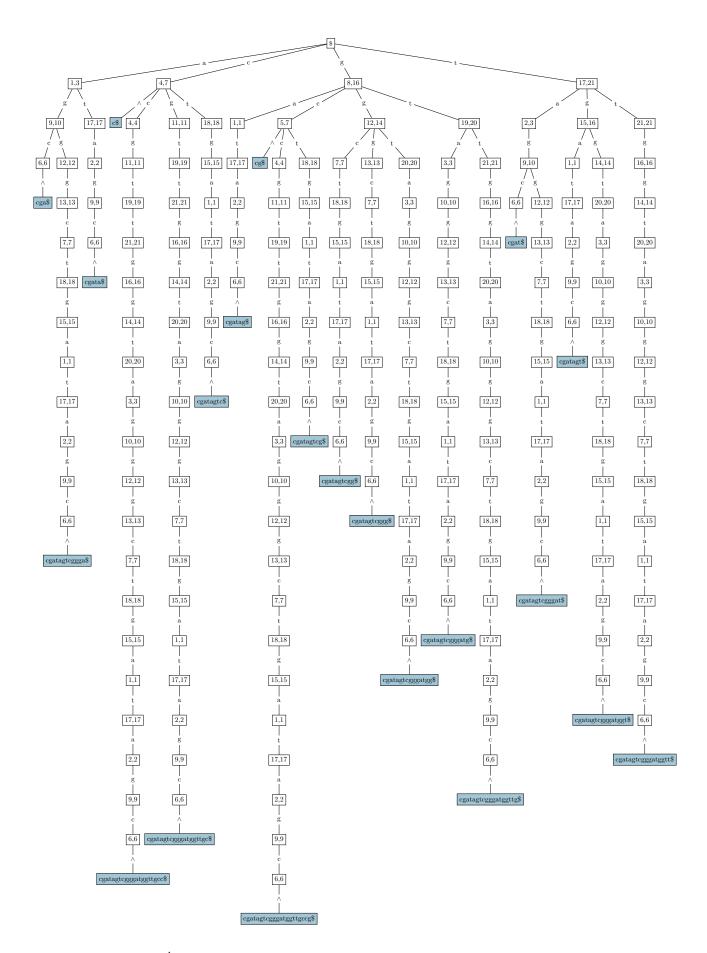


Figura 4: Árbol de prefijos con intervalos SA de la palabra "cgatagtcgggatggttgccg\$".

3. Anejo

Listing 1: python inexrecur_clean.py

```
#!/usr/bin/env python
2
   # Started Wed 25 Mar 00:56:18 2020 by nahnero. #
   def CalculateD (W):
3
4
       D = [0]*len (W)
       z, j = 0, 0
5
6
       for i in range (len (W)):
7
            if W[j:i] not in X:
8
                z += 1
                j = i + 1
9
10
            D[i] = z
11
       return D
12
   def InexRecur (W, i, z, k, 1):
13
       if i < 0:
14
            if z < 0:
15
16
                return set ([])
17
            else:
                return set ([(k, 1)])
18
19
       if z < D[i]:
20
           return set ([])
21
       I = set ([])
22
       I = I.union (InexRecur (W, i-1, z-1, k, 1))
23
       for bi in dic:
24
            ki = C (bi) + O (bi, k-1) + 1
25
            li = C (bi) + O (bi, 1)
26
            if ki <= li:
27
                I = set (I.union (InexRecur (W, i, z-1, ki, li)))
28
                if bi == W[i]:
29
                    I = set (I.union (InexRecur (W, i-1, z, ki, li)))
30
                else:
                    I = set (I.union (InexRecur (W, i-1, z-1, ki, li)))
31
32
       return I
33
34
   def InexSearch (W, z):
35
       return InexRecur (
36
                W = W,
37
                i = len (W)-1,
38
                z = z,
39
                k = 0,
40
                l = len (X)-1)
41
   X = 'googol$';
   W = 'lol'
42
43
   dic = ''.join (sorted (set (X.replace ('$', ''))))
44
   BWT = [[(X*2)[i:i+len (X)], i] for i in range (len (X))]
45
   BWT.sort ()
46
47
   S, B = [], ''
48
   for i in BWT: S.append (i[1]); B += (i[0][-1])
   C = lambda \ a: \ ''.join (sorted (X.replace ('$', '')))[0:-2].find (a)
49
50
   0 = lambda a, i: B[0:i+1].count (a)
51
52
   D = CalculateD (W)
53
54
   IS = InexSearch (W, 1)
55
   for i in IS: print (i)
```

```
1
   #!/usr/bin/env Rscript
 2
   C <- function(a) {</pre>
 3
      components <- gsub("[[:punct:]]","",X)</pre>
      components <- strsplit(components,"")</pre>
 4
      components_sorted <- sort(components[[1]])</pre>
 5
 6
      idx <- min(which(components_sorted == a))</pre>
 7
      C_num < - idx - 1
 8
      return(C_num)
9
10
   0 <- function(a,i) {</pre>
11
      i < -i + 1
12
      coincidences <- grep(a,B)</pre>
      counter <- i>= coincidences
13
14
      O_num <- length(counter[counter==TRUE])</pre>
15
      return(O_num)
16
   }
17
    INEXRECUR <- function(W,i,z,k,l,sep) {</pre>
18
      if (i<0){
19
20
        if (z<0) {return({})}
21
        else {
22
          result <- list(c(k,1))
23
          return(result)
24
        }}
25
      if (z<D[i+1]) {return({})}
26
      I <- {}
27
      I <- c (I, INEXRECUR(W,i-1,z-1,k,l,sep));</pre>
28
      for (bi in b){
29
        ki \leftarrow C(bi) + O(bi, k-1) + 1 #Precalculate C and O function
        li <- C(bi) + O(bi,l) #Precalculate C and O function
30
31
        if (ki<=li){
32
          I \leftarrow c (I, INEXRECUR(W,i,z-1,ki,li,sep));
33
          if (bi == W[[1]][i+1]){
34
             I <- c (I, INEXRECUR(W,i-1,z,ki,li,sep))</pre>
35
             I <- c (I, INEXRECUR(W,i-1,z-1,ki,li,sep))</pre>
36
37
          }}}
      return(I)
38
39
   } #end INEXRECUR
40
41
42
   X <- "googol$"
   SA_pre <- c()
43
   SA_pre[1] <- X
44
   for (i in 2:nchar(X)){
45
46
   SA_pre[i] <- paste(substring(X,i,nchar(X)),</pre>
                         substring(X,1,i-1),sep='',collapse=NULL)
47
48
49
   SA_list <-sort(SA_pre, index.return=TRUE)</pre>
   SA<-SA_list$x
50
   S<-SA_list$ix
51
   B<-paste(substring(SA,nchar(X),nchar(X)),sep='',collapse=NULL)
52
   b <- c("g","1","o")
53
54
   D \leftarrow c(0,0,1)
55
   W = list(c("l", "o", "l"))
56
   i = 2
57
   z = 1
   a = INEXRECUR(W,i,z,0,6,"")
58
59
   print (a)
```

```
1
   int main (void){
            char X[] = "googol$";
            char W[] = "lol";
3
            char dict[] = "glo";
4
5
            setup (X, W, dict);
6
            return 0;
7
8
9
10
   void setup (char* X1, char* W1, char* dic1){
11
            int i, length; i = length = 0;
12
13
                  = X1;
14
            /* reverse */
15
                  = W1;
16
            while (W1[++i]);
17
            Wlen = i;
18
            dic
                  = dic1;
19
            while (X[++length]);
20
21
22
                = (int*) calloc (length, sizeof (int));
23
            BWT = (char**)calloc (length, sizeof (char*));
24
               = (char*) calloc (length, sizeof (char));
25
                = (char*) calloc (length, sizeof (char));
26
            X2 = (char*) calloc (2*length+1, sizeof (char));
            for (i = 0; i != length; i++)
27
28
                    BWT[i] = (char*)malloc (length*sizeof (char)+1);
29
30
           sprintf (X2, "%s%s", X, X);
31
32
            for (i = 0; i != length; i++)
33
                    sprintf (BWT[i], "%.*s", (int)length, X2+i);
34
35
36
            sort (BWT, length);
37
            CalculateD (W1);
38
39
            for (i = 0; i != length; i++) B [i] = BWT[i][length-1];
40
            for (i = 0; i != length; i++) Xs[i] = BWT[i][0];
41
42
            i = 0; while (X[++i]);
            InexRecur (W1, Wlen - 1, 1, 0, i - 1);
43
44
45
            free (Xs);
46
            free (X2);
47
            free (D);
48
            for (i = 0; i != length ; i++) free (BWT[i]);
49
            free (BWT);
50
51
52
   static int myCompare(const void* a, const void* b){
53
           return strcmp (*(const char**)a, *(const char**)b);
54
   }
55
56
   void sort(char* arr[], int n){
57
            qsort (arr, n, sizeof(char*), myCompare);
   }
58
```

```
59
60
    int C (char a){
             int i = 0;
61
             while (X[++i]) if (Xs[i] == a) return i - 1;
62
63
             return 0;
64
65
66
    int 0 (char a, int i){
67
             int acc = 0, j;
             for (j = 0; j \le i; j++) {
68
                     if (B[j] == a) acc++;
69
70
71
             return acc;
72
    }
73
74
    /* Ignacio Amat */
75
76
    void CalculateD (char* W){
77
             short z = 0, j = 0, i;
78
             char buf[BUFSIZ] = \{'\0'\};
79
             extern int Wlen;
80
             for (i = 0; i != Wlen; i++){}
81
                      sprintf (buf, "%.*s", i, W+j);
82
                      if (!strstr (X, buf)){
83
                              z++;
                              j = i + 1;
84
85
                     D[i] = z;
86
87
             }
88
    }
89
90
    void InexRecur (char* W, int i, int z, int k, int l){
91
             if (i < 0){
92
                      if (z < 0){
93
                              return;
94
                      } else {
95
                              printf ("(%d, %d)\n", k, 1);
96
                              return;
97
                     }
             }
98
             if (z < D[i]){return;}</pre>
99
100
             InexRecur (W, i-1, z-1, k, 1);
101
             int ki, li, m;
102
             for (m = 0; m != 3; m++){
103
                     ki = C (dic [m]) + O (dic [m], k-1) + 1;
104
                     li = C (dic [m]) + O (dic [m], 1);
105
                      if (ki <= li){
106
                               InexRecur (W, i, z-1, ki, li);
                              if (dic [m] == W[i]){
107
108
                                       InexRecur (W, i-1, z, ki, li);
109
                              } else {
110
                                       InexRecur (W, i-1, z-1, ki, li);
111
                              }
112
                     }
113
             }
114
             return;
115
```