

## CSE 5311 Notes 15: Sequences

(Last updated 8/10/15 3:14 PM)

### PATTERN-BASED PREPROCESSING

#### Simple Rescanning

Pattern -  $m$  symbols    Text -  $n$  symbols

Example:    Pattern:    0 1 2 3  
                           A B A C  
                           Text:    A B A B A B A C A B A B A C A B A  
                                   0 1 2 3  
                                   0  
                                   0 1 2 3  
                                   0  
                                   0 1 2 3  
                                   0  
                                   0 1  
                                   0  
                                   0 1 2 3  
                                   0  
                                   0 1 2 3  
                                   0  
                                   0 1

Worst-Case:  $\Theta(mn)$

                          0 1 2 3  
                           Pattern: 0 0 0 1  
                           Text:    0 0 0 0 0 0 0 0 0 0 1  
                                   0 1 2 3  
                                   0 1 2 3  
                                   0 1 2 3  
                                   0 1 2 3  
                                   0 1 2 3  
                                   0 1 2 3  
                                   0 1 2 3  
                                   0 1 2 3

#### Rabin-Karp Algorithm

Concept:

1. Compute signature of all  $m$  symbols of pattern. (Taking  $\theta(m)$  time.)
2. Compute signature of each  $m$  contiguous symbols of text. (Taking overall  $\theta(n)$  time.)

If function values for 1. and 2. are ever equal, then do symbol-by-symbol test.

```

#include <stdio.h>
#include <string.h>

#define B 131
#define EOS 0
#define TRUE 1
#define FALSE 0

void search(pat,text)
char *pat,*text;
{
    int hpat,htext,Bm,j,m;

    if (pat[0]==EOS) return;
    Bm=1;
    hpat=htext=0;

    for (m=0;text[m]!=EOS && pat[m]!=EOS;m++)
    {
        Bm*=B;
        hpat=hpat*B+pat[m];
        htext=htext*B+text[m];
        printf("%c ",text[m]);
        if (pat[m+1]==EOS)
            printf("%d",htext);
        else
            printf("\n");
    }

    if (text[m]==EOS && pat[m]!=EOS) return;

    for (j=m;TRUE;j++)
    {
        if (hpat==htext && strncmp(text+j-m,pat,m)==0)
            printf("<<<\n");
        else
            printf("\n");
        if (text[j]==EOS) return;
        htext=htext*B-text[j-m]*Bm+text[j];
        printf("%c %d ",text[j],htext);
    }
}

main()
{
    char pat[80],text[80],*pos,*work;

    printf("Enter pattern & text\n");
    while (scanf("%s %s",pat,text)!=EOF)
    {
        search(pat,text);
        printf("Enter pattern & text\n");
    }
}

```



## Knuth-Morris-Pratt Failure Link Construction

Fail link - seeks to reuse largest possible *suffix* before present position that matches a *prefix* of pattern.

Style 1: Choose maximum value of  $k$  with  $0 \leq k < j$  such that

for  $0 < i \leq k$ :  $\text{pattern}[k - i] == \text{pattern}[j - i]$

Now set  $\text{fail}[j] = k$

Style 2: Choose maximum value of  $k$  with  $0 \leq k < j$  such that

$\text{pattern}[k] != \text{pattern}[j]$ , and

for  $0 < i \leq k$ :  $\text{pattern}[k - i] == \text{pattern}[j - i]$

Now set  $\text{fail}[j] = k$

Direct application of these definitions could take  $\Theta(m^3)$  time!

Either style fail link table may be constructed in  $\Theta(m)$  time.

For style 1:

Suppose fail links 0 through  $j$  have already been set and  $\text{fail}[j] == k$ .

If, in addition,  $\text{pattern}[j] == \text{pattern}[k]$  then

Set  $\text{fail}[j+1] = k+1$

But, what if  $\text{pattern}[j] != \text{pattern}[k]$ ?

Move  $k$  back to  $\text{fail}[k]$  and recheck for pattern match

0 b -1	11 b 6
1 a 0	12 a 4
2 b 0	13 b 5
3 b 1	14 a 6
4 a 1	15 b 7
5 b 2	16 b 8
6 a 3	17 a 9
7 b 2	18 b 10
8 b 3	19 a 11
9 a 4	20 b ?
10 b 5	

For style 2:

Suppose fail links 0 through  $j$  have already been set and  $k$  is the *required maximum value* that was used when setting  $\text{fail}[j]$ .

If, in addition,  $\text{pattern}[j] == \text{pattern}[k]$  then

If  $\text{pattern}[j+1] != \text{pattern}[k+1]$

Set  $\text{fail}[j+1] = k+1$

Else

Set  $\text{fail}[j+1] = \text{fail}[k+1]$

But, what if  $\text{pattern}[j] != \text{pattern}[k]$ ?

Move  $k$  back to  $\text{fail}[k]$  and recheck for pattern match

0 b -1	7 b -1	14 a 3
1 a 0	8 b 1	15 b -1
2 b -1	9 a 0	16 b 1
3 b 1	10 b -1	17 a 0
4 a 0	11 b 6	18 b -1
5 b -1	12 a 0	19 a 11
6 a 3	13 b -1	20 b ?

	Fail 1	Pattern	Fail 2		Fail 1	Pattern	Fail 2
0		A		0		A	
1		A		1		B	
2		B		2		A	
3		A		3		A	
4		A		4		B	
5		B		5		A	
6		A		6		B	
7		A		7		A	
8		A		8		A	
9		B		9		B	
				10		A	
				11		A	
				12		B	

## KMP.c:

```

/* Determine all possible occurrences of pattern string in text
string using KMP technique*/
#include <stdio.h>
#include <string.h>

#define MAXPATLEN 80
#define EOS 0

void preprocpat1(pat,next)
/* Produces slower failure links, but capable of continuing after
match */
char *pat;
int next[];
{
int i,j;
i=0;
j=next[0]= -1;
do
{
if (j==(-1)||pat[i]==pat[j])
{
i++;
j++;
next[i]=j;
}
else
j=next[j];
} while (pat[i]!=EOS);
printf("Fail link table 1\n");
for (i=0;pat[i]!=EOS;i++)
printf("%d %c %d\n",i,pat[i],next[i]);
}

void preprocpat2(pat,next)
/* Produces faster failure links, but INCAPABLE of continuing after
match */
char *pat;
int next[];
{
int i,j;
i=0;
j=next[0]= -1;
do
{
if (j==(-1)||pat[i]==pat[j])
{
i++;
j++;
next[i]=(pat[j]==pat[i]) ? next[j] : j;
}
else
j=next[j];
} while (pat[i]!=EOS);
printf("Fail link table 2\n");
for (i=0;pat[i]!=EOS;i++)
printf("%d %c %d\n",i,pat[i],next[i]);
}

void match(text1,pat)
char *text1,*pat;
{
int next1[MAXPATLEN],next2[MAXPATLEN],i,j;
char *text;
char printSymbol=' ';

if (*pat==EOS) return;
preprocpat1(pat,next1);
preprocpat2(pat,next2);

printf("%s\n",text1);
text=text1;

```

```

for (j=0; *text!=EOS;)
  if (j==(-1) || pat[j]== *text)
  {
    if (pat[j+1]==EOS)
      printSymbol='^';
    if (pat[j+1]==EOS && *(text+1)!=EOS)
      j=next1[j]; /*restart*/
    else
    {
      printf("%c",printSymbol);
      printSymbol=' ';
      text++;
      j++;
    }
  }
  else
    j=next2[j];
printf("\n");
}

main()
{
char pat[80],text[80];

printf("Enter text & pattern\n");
/*variable text is string 1, variable pat is string 2*/
while (scanf("%s %s",text,pat)!=EOF)
{
  match(text,pat);
  printf("Enter text & pattern\n");
}
}

```

## KMP Complexity

$m$  pattern symbols,  $n$  text symbols

Matcher - every comparison is preceded by a movement of one or both pointers

Text pointer always moves forward  $\Rightarrow \Theta(n)$

Pattern pointer

Moves forward once for each text symbol  $\Rightarrow \Theta(n)$

Total number of backward movements  $\leq$  total number of forward movements  
 $\Rightarrow O(n)$  (could set up potential function based on pattern position)

Failure table construction

Lead pointer always moves forward  $\Rightarrow \Theta(m)$

Prefix pointer

Moves forward once for each pattern symbol  $\Rightarrow \Theta(m)$

Total number of backward movements  $\leq$  total number of forward movements  
 $\Rightarrow O(m)$  (could set up potential function based on pattern position)

Overall:  $\Theta(m + n)$

Aside: Worst-case number of comparisons (fail 2 links) when processing a text symbol is bounded by  $1 + 1.44 \lg m$ .

Fibonacci strings may be used as difficult cases:

$$F_1 = \text{"a"} \qquad F_2 = \text{"b"} \qquad F_n = F_{n-1}F_{n-2}$$

$$F_3 = \text{"ba"}$$

$$F_4 = \text{"bab"}$$

$$F_5 = \text{"babba"}$$

$$F_6 = \text{"babbabab"}$$

$$F_7 = \text{"babbababbabba"}$$

#### Other Applications

Is  $W$  constructed by repetition? i.e.  $W = X^k$  where  $X$  is the *smallest* such string.

String length must be factorable, then use KMP (either method) and use pointers at end.

	Pattern	Fail 2		Pattern	Fail 2
0	0	-1	0	1	-1
1	0	-1	1	2	0
2	1	1	2	3	0
3	0	-1	3	4	0
4	0	-1	4	1	-1
5	2	2	5	2	0
6	0	-1	6	3	0
7	0	-1	7	4	0
8	1	1	8	1	-1
9	0	-1	9	2	0
10	0	-1			
11	2	2			

$$\frac{12}{12-6} = 2$$

$$\frac{10}{10-6} = 2.5$$



Are two strings (of same length) “circularly equal”?

1. Text = string2 || string 2
2. Pattern = string1
3. Does pattern occur in text?

Aside: CLRS 32.4 Prefix Function

Use Pascal, not C, array conventions:

	Pattern	Fail 1	Fail 2	$\pi$
1	a	0	0	0
2	b	1	1	0
3	a	1	0	1
4	b	2	1	2
5	a	3	0	3
6	b	4	1	4
7	a	5	0	5
8	b	6	1	6
9	c	7	7	0
10	a	1	0	1
11	?	2		

$$\pi[i] = \text{Fail1}[i+1] - 1$$

Gusfield’s Z Algorithm (“Fundamental String Preprocessing”)

Alternative to KMP

More straightforward to apply for new situations.

Instead of failure links, each position of string has the number of positions starting at that point that match a prefix of entire string.

Trivial to construct Z table in  $\Theta(n^2)$  time, but can do in  $\Theta(n)$  time.

	Pattern	Z
0	0	8
1	1	0
2	0	1
3	0	3
4	1	0
5	0	3
6	1	0
7	0	1

Applications:

String search: Build Z table for <pattern><text>

Any symbol in <text> with a Z value  $\geq$  |<pattern>| gives a match.

Circularly equal: Build Z table for <string1><string2><string2> (|<string1>| = |<string2>|)

Any symbol in first copy of <string2> with a Z value  $\geq$  |<string2>| verifies.

Three aspects to Z-table construction:

1. Scanning to match against prefix.
2. Testing whether scanning must be resumed or prefix Z value copying may continue (“Z box” and “countdown”).
3. Modified copying (second phase) when scanning has reached end of input string.

0	1	2	3	4	5	6	7	8	9	10	11	12
0	1	0	0	1	0	1	0	0	1	0	0	1

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
a	b	c	a	b	a	b	c	a	b	c	a	b	a	b	c	a	b	a	b	c

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
0	1	0	0	1	0	1	0	0	1	0	0	1	0	1	0	0	1	0	1	0

21	22	23	24	25	26	27	28	29	30	31	32	33
0	1	0	0	1	0	1	0	0	1	0	0	1

```

void fundamental() {
int fillPosition,scanPosition,countdown,prefixPosition;

sLength=strlen(S);
if (sLength==0)
    return;
Z[0]=sLength;

//Before EOS scanned
countdown=0;
fillPosition=scanPosition=1;
while (scanPosition<sLength)
    if (countdown==0 || Z[prefixPosition]>=countdown) {
        while (S[scanPosition]==S[countdown]) {
            scanPosition++;
            countdown++;
        }
        Z[fillPosition++]=countdown;
        if (countdown==0)
            scanPosition=fillPosition;
        else {
            countdown--;
            prefixPosition=1;
        }
    }
    else {
        Z[fillPosition++]=Z[prefixPosition++];
        countdown--;
    }

//After EOS scanned
while (fillPosition<sLength)
    Z[fillPosition++]=min(countdown--,Z[prefixPosition++]);
}

```

S. Faro & T. Lecroq, “The Exact Online String Matching Problem: A Review of the Most Recent Results”, *ACM Computing Surveys* 45 (2), Feb 2013:  
<http://dl.acm.org.ezproxy.uta.edu/citation.cfm?doid=2431211.2431212>

## SUFFIX ARRAYS: TEXT-BASED PREPROCESSING

Concept: Preprocess text and handle different patterns on-the-fly.

Simpler and more compact than other structures, such as suffix trees.

Sort array (**sa**) of pointers/subscripts based on corresponding suffixes of a text.

	0	1	2	3	4	5	6	7	8	9	10	11
s	a	b	c	d	a	b	c	d	a	b	c	\0
sa	11	8	4	0	9	5	1	10	6	2	7	3

(These notes use null-terminated strings. Other terminators, e.g. \$, would change the results . . .)

Key-comparison sorts can construct in  $O(n^2 \log n)$  time (expected time for qsort):

```
int suffixCompare(const void *xVoidPt, const void *yVoidPt) {
// Used in qsort call to generate suffix array.
int *xPt=(int*) xVoidPt, *yPt=(int*) yVoidPt;

return strcmp(&s[*xPt], &s[*yPt]);
}

scanf("%s", s);
n=strlen(s)+1;
for (i=0; i<n; i++)
    sa[i]=i;
qsort(sa, n, sizeof(int), suffixCompare);
```

The *rank* of a suffix is its position in the suffix array (i.e. these are inverses:  $\text{rank}[\text{sa}[i]] = i$ ).

The *longest common prefix* ( $\text{lcp}[i]$ ) is the number of prefix matches for  $\text{sa}[i-1]$  and  $\text{sa}[i]$ .

i	sa	suffix	lcp	s	rank	lcp[rank]
0	11		-1	a	3	7
1	8	abc	0	b	6	6
2	4	abcdabc	3	c	9	5
3	0	abcdabcdabc	7	d	11	4
4	9	bc	0	a	2	3
5	5	bcdabc	2	b	5	2
6	1	bcdabcdabc	6	c	8	1
7	10	c	0	d	10	0
8	6	cdabc	1	a	1	0
9	2	cdabcdabc	5	b	4	0
10	7	dabc	0	c	7	0
11	3	dabcdabc	4		0	-1

U. Manber and E. Myers, "Suffix Arrays: A New Method for On-Line String Searches", *SIAM J. on Computing* 22, 5 (1993), 935-948. <http://epubs.siam.org.ezproxy.uta.edu/doi/pdf/10.1137/0222058>

J. Kärkkäinen et al., "Linear Work Suffix Array Construction", *J. ACM* 53, 6 (Nov. 2006), 918-936 achieves linear time. <http://dl.acm.org.ezproxy.uta.edu/citation.cfm?doid=1217856.1217858>

S.J. Puglisi et al., "A Taxonomy of Suffix Array Construction Algorithms", *ACM Computing Surveys* 39 (2), June 2007. <http://dl.acm.org.ezproxy.uta.edu/citation.cfm?doid=1242471.1242472>

*"Impressive as the progress has been, ingenious as the methods have been, there still remains the challenge to devise a SACA that is lightweight, linear in the worst case, and fast in practice."*

T. Kasai et al., "Linear-Time Longest-Common-Prefix Computation in Suffix Arrays and Its Applications", p. 181-192 in A. Amir and G.M. Landau (Eds): *CPM 2001: 12th Ann'l Symp. on Combinatorial Pattern Matching*, Lecture Notes in Computer Science 2089, Springer-Verlag, 2001.

M.I. Abouelholda, et al., "Replacing Suffix Trees with Enhanced Suffix Arrays", *J. of Discrete Algorithms* 2 (2004), 53-86. <http://www.sciencedirect.com.ezproxy.uta.edu/science/journal/15708667/2/1>

G. Navarro, "Spaces, Trees, and Colors: The Algorithmic Landscape of Document Retrieval on Sequences", *ACM Computing Surveys* 46 (4), March 2014. <http://dl.acm.org.ezproxy.uta.edu/citation.cfm?doid=2597757.2535933>

Manber/Myers method uses  $O(\log n)$  radix sorts each taking  $O(n)$  time to achieve  $O(n \log n)$  time.

radixsort i=1			POS			MSD LSD - after radix sort		
POS	MSD	LSD - start of pass	11			1	0	
0 abcdabcdabc	98	99	0 abcdabcdabc			98	99	
1 bcdabcdabc	99	100	4 abcdabc			98	99	
2 cdabcdabc	100	101	8 abc			98	99	
3 dabcdabc	101	98	1 bcdabcdabc			99	100	
4 abcdabc	98	99	5 bcdabc			99	100	
5 bcdabc	99	100	9 bc			99	100	
6 cdabc	100	101	10 c			100	1	
7 dabc	101	98	2 cdabcdabc			100	101	
8 abc	98	99	6 cdabc			100	101	
9 bc	99	100	3 dabcdabc			101	98	
10 c	100	1	7 dabc			101	98	
11	1	0	6 buckets					
POS			radixsort i=2			MSD LSD - start of pass		
11	MSD	LSD - after renumbering	POS			MSD	LSD - start of pass	
0 abcdabcdabc	1	0	0 abcdabcdabc			1	4	
4 abcdabc	1		1 bcdabcdabc			2	5	
8 abc	1		2 cdabcdabc			4	1	
1 bcdabcdabc	2		3 dabcdabc			5	2	
5 bcdabc	2		4 abcdabc			1	4	
9 bc	2		5 bcdabc			2	5	
10 c	3		6 cdabc			4	1	
2 cdabcdabc	4		7 dabc			5	2	
6 cdabc	4		8 abc			1	3	
3 dabcdabc	5		9 bc			2	0	
7 dabc	5		10 c			3	0	
			11			0	0	
POS			POS			MSD LSD - after renumbering		
11	MSD	LSD - after radix sort	11			MSD	LSD - after renumbering	
8 abc	1	3	8 abc			1		
0 abcdabcdabc	1	4	0 abcdabcdabc			2		
4 abcdabc	1	4	4 abcdabc			2		
9 bc	2	0	9 bc			3		
1 bcdabcdabc	2	5	1 bcdabcdabc			4		
5 bcdabc	2	5	5 bcdabc			4		
10 c	3	0	10 c			5		
2 cdabcdabc	4	1	2 cdabcdabc			6		
6 cdabc	4	1	6 cdabc			6		
3 dabcdabc	5	2	3 dabcdabc			7		
7 dabc	5	2	7 dabc			7		
8 buckets								
radixsort i=4			POS			MSD LSD - after radix sort		
POS	MSD	LSD - start of pass	11			MSD	LSD - after radix sort	
0 abcdabcdabc	2	2	8 abc			0	0	
1 bcdabcdabc	4	4	4 abcdabc			1	0	
2 cdabcdabc	6	6	0 abcdabcdabc			2	1	
3 dabcdabc	7	7	9 bc			2	2	
4 abcdabc	2	1	5 bcdabc			3	0	
5 bcdabc	4	3	1 bcdabcdabc			4	3	
6 cdabc	6	5	10 c			4	4	
7 dabc	7	0	6 cdabc			5	0	
8 abc	1	0	2 cdabcdabc			6	5	
9 bc	3	0	7 dabc			6	6	
10 c	5	0	3 dabcdabc			7	0	
11	0	0	12 buckets			7	7	
POS			POS			MSD LSD - after renumbering		
11	MSD	LSD - after renumbering	11			MSD	LSD - after renumbering	
8 abc	1		8 abc			1		
4 abcdabc	2		4 abcdabc			2		
0 abcdabcdabc	3		0 abcdabcdabc			3		
9 bc	4		9 bc			4		
5 bcdabc	5		5 bcdabc			5		
1 bcdabcdabc	6		1 bcdabcdabc			6		
10 c	7		10 c			7		
6 cdabc	8		6 cdabc			8		
2 cdabcdabc	9		2 cdabcdabc			9		
7 dabc	10		7 dabc			10		
3 dabcdabc	11		3 dabcdabc			11		

The lcp array may be constructed in  $O(n)$  time. (LCPdemo.c)

```

void computeLCP() {
//Kasai et al linear-time construction
int h,i,j,k;

h=0; // Index to support result that lcp[rank[i]]>=lcp[rank[i-1]]-1
for (i=0;i<n;i++) {
    k=rank[i];
    if (k==0)
        lcp[k]=(-1);
    else {
        j=sa[k-1];
        // Attempt to extend lcp
        while (i+h<n && j+h<n && s[i+h]==s[j+h])
            h++;
        lcp[k]=h;
    }
    if (h>0)
        h--; // Decrease according to result
}
}

```

Accelerating binary searches on sa using lcp:

```

int slowSearchFirst(char* key) {
// Finds first string >= key assisted by binary search on suffix array,
// but without using lcp.
int low,high,mid;
int i,j;

low=0;
high=n-1;
while (low<=high) {
    mid=(low+high)/2;
    j=sa[mid]; // Position in s
    i=0; // Position in key
    // Like strcmp
    while (s[j]==key[i] && key[i]) {
        j++;
        i++;
    }
    if (key[i]==0 || s[j]>key[i])
        high=mid-1;
    else
        low=mid+1;
}

return low;
}

```

```

int fastSearchFirst(char* key) {
// Finds first string >= key assisted by binary search on suffix array.
// Uses lcp to limit char comparisons.
int low,high,mid;
int i,j;
int keyLength,lowMatches,highMatches,midMatches;

keyLength=strlen(key);
low=0;
lowMatches=0;
high=n-1;
highMatches=0;

while (low<=high) {
mid=(low+high)/2;
// How many strcmp matches can be salvaged?
midMatches=(lowMatches<highMatches) ? lowMatches : highMatches;

if (midMatches==keyLength) {
high=mid-1;
highMatches=(lcp[mid]<midMatches) ? lcp[mid] : midMatches;
}

j=sa[mid]+midMatches; // Position in s
i=midMatches; // Position in key
while (s[j]==key[i] && key[i]) {
midMatches++;
j++;
i++;
}
if (key[i]==0 || s[j]>key[i]) {
high=mid-1;
highMatches=(lcp[mid]<midMatches) ? lcp[mid] : midMatches;
}
else {
low=mid+1;
lowMatches=(lcp[low]<midMatches) ? lcp[low] : midMatches;
}
}

return low;
}

```

slowSearchFirst uses  $\Theta(|key| \cdot \log|text|)$  time, fastSearchFirst uses  $\Theta(|key| + \log|text|)$  time.

### Longest Common Substring Using Suffix Array & LCP

1. Construct text as  $\langle \text{string1} \rangle \$ \langle \text{string2} \rangle$ , then build SA and LCP arrays.
2. Scan LCP array for maximum entry for two adjacent SA entries from different strings.

i	sa	suffix	lcp	s	rank	lcp[rank]
0	30		-1	0	13	8
1	11	\$010010100101001001	0	1	24	7
2	10	0\$010010100101001001	0	0	5	6
3	27	001	1	0	14	6
4	24	001001	3	1	25	5
5	2	001001010\$010010100101001001	6	0	6	4
6	5	001010\$010010100101001001	4	0	17	3
7	19	00101001001	6	1	28	2
8	14	0010100101001001	9	0	10	2
9	28	01	1	1	21	1
10	8	010\$010010100101001001	2	0	2	0
11	25	01001	3	\$	1	0
12	22	01001001	5	0	16	11
13	0	01001001010\$010010100101001001	8	1	27	10
14	3	01001010\$010010100101001001	6	0	8	9
15	17	0100101001001	8	0	19	8
16	12	010010100101001001	11	1	30	7
17	6	01010\$010010100101001001	3	0	15	8
18	20	0101001001	5	1	26	7
19	15	010100101001001	8	0	7	6
20	29	1	0	0	18	5
21	9	10\$010010100101001001	1	1	29	4
22	26	1001	2	0	12	5
23	23	1001001	4	1	23	4
24	1	1001001010\$010010100101001001	7	0	4	3
25	4	1001010\$010010100101001001	5	0	11	3
26	18	100101001001	7	1	22	2
27	13	10010100101001001	10	0	3	1
28	7	1010\$010010100101001001	2	0	9	1
29	21	101001001	4	1	20	0
30	16	10100101001001	7	0	-1	

Length of longest common substring is 8  
01001001

M.A. Babenko and T.A. Starikovskaya, "Computing Longest Common Substrings Via Suffix Arrays", *Proc. 3rd Int'l Computer Science Symp. in Russia, CSR 2008, LNCS 5010, Springer, 64-75.*

Aside: Kärkkäinen et.al. linear-time suffix array construction (DC3)

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21  
c d e a b c d e c d e a b c d e a b c d e (called t in the paper)

**Step 1: Sort Sample Suffixes**

$R_1 = [dea][bcd][ecd][eab][cde][abc][de0]$  (Sample triples starting at position 1)

$R_2 = [eab][cde][cde][abc][dea][bcd][e00]$  (Sample triples starting at position 2)

Concatenate to give the "samples":

$R = [dea][bcd][ecd][eab][cde][abc][de0][eab][cde][cde][abc][dea][bcd][e00]$



LSD Radix sort the “characters” of R and rank:

```
[000] 1
[abc] 2 (appears 2 times)
[bcd] 3 (appears 2 times)
[cde] 4 (appears 3 times)
[de0] 5
[dea] 6 (appears 2 times)
[e00] 7
[eab] 8 (appears 2 times)
[ecd] 9
```

R' is the result of substituting the rank for each character in R:

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14
6 3 9 8 4 2 5 1 8 4 4 2 6 3 7
```

Since the symbols in R' are not unique, DC3 is called recursively to produce the suffix array  $SA_{R'}$ . Note that 15 is the empty suffix.

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14
7 5 11 13 1 4 10 9 6 12 0 14 3 8 2
```

From Kärkkäinen et.al., “Once the sample suffixes are sorted, assign a rank to each suffix. . .”. Each value appearing in  $SA_{R'}$  gets mapped as below according to how R was constructed from the triples:

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 (Positions of sample)
1 4 7 10 13 16 19 22 2 5 8 11 14 17 20 (Positions of original Fibonacci string)
```

to give the following when substituted into  $SA_{R'}$ :

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14
7 5 11 13 1 4 10 9 6 12 0 14 3 8 2 (Positions of sample)
22 16 11 17 4 13 8 5 19 14 1 20 10 2 7 (Positions of original Fibonacci string)
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 (Ranks)
```

Now “invert” to give ranks for sampled positions of original Fibonacci string:

```
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
c d e a b c d e c d e a b c d e a b c d e
? 11 14 ? 5 8 ? 15 7 ? 13 3 ? 6 10 ? 2 4 ? 9 12 ? 1
```

### Step 2: Sort Nonsample Suffixes

Nonsample suffixes are indicated in the immediately preceding table by “?”. Use LSD radix sort to sort pairs  $(t_i, \text{rank}(S_{i+1}))$  for these positions:

Sort input:

```
0: (c, 11)
3: (a, 5)
6: (d, 15)
9: (d, 13)
12: (b, 6)
15: (e, 2)
18: (c, 9)
21: (0, 1)
```

Sort output

- 21: (0, 1)
- 3: (a, 5)
- 12: (b, 6)
- 18: (c, 9)
- 0: (c, 11)
- 9: (d, 13)
- 6: (d, 15)
- 15: (e, 2)

**Step 3: Merge**

- |               |            |            |
|---------------|------------|------------|
| 16: (a, 4)    | 21: (0, 1) | (0, 0, 0)  |
| 11: (a, b, 6) | 3: (a, 5)  | (a, b, 8)  |
| 17: (b, c, 9) | 12: (b, 6) | (b, c, 10) |
| 4: (b, 8)     | 18: (c, 9) | (c, d, 12) |
| 13: (c, 10)   | 0: (c, 11) | (c, d, 14) |
| 8: (c, d, 13) | 9: (d, 13) | (d, e, 3)  |
| 5: (c, d, 15) | 6: (d, 15) | (d, e, 7)  |
| 19: (d, 12)   | 15: (e, 2) | (e, a, 4)  |
| 14: (d, e, 2) |            |            |
| 1: (d, 14)    |            |            |
| 20: (e, 0, 1) |            |            |
| 10: (e, 3)    |            |            |
| 2: (e, a, 5)  |            |            |
| 7: (e, 7)     |            |            |

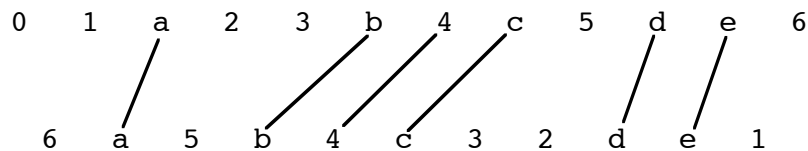
**Final Result**

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	
21	16	11	3	17	12	4	18	13	8	0	5	19	14	9	1	6	20	15	10	2	7	Suffix Array
-1	0	5	5	0	4	4	0	3	8	8	3	0	2	7	7	2	0	1	6	6	1	LCP

LONGEST COMMON SUBSEQUENCES

Dynamic Programming - review

Has important applications in genetics research.



1. Describe problem input.

Two sequences:

$$X = x_1 x_2 \dots x_m$$

$$Y = y_1 y_2 \dots y_n$$

2. Determine cost function and base case.

$C(i, j)$  = length of LCS for  $x_1 x_2 \dots x_i$  and  $y_1 y_2 \dots y_j$

$C(i, j) = 0$  if  $i = 0$  or  $j = 0$

3. Determine general case.

Suppose  $C(i, j)$  has

$$x_1 x_2 \dots x_{i-1}^A \quad y_1 y_2 \dots y_{j-1}^A$$

Since  $x_i = y_j$ ,  $C(i, j) = C(i-1, j-1) + 1$

Now suppose  $x_i \neq y_j$ :

$$x_1 x_2 \dots x_{i-1}^A \quad y_1 y_2 \dots y_{j-1}^B$$

But 'B' may appear in  $x_1 x_2 \dots x_{i-1}$  or 'A' may appear in  $y_1 y_2 \dots y_{j-1}$ :

$$C(i, j) = \max\{C(i, j-1), C(i-1, j)\} \text{ if } x_i \neq y_j$$

4. Appropriate ordering of subproblems.

Before computing  $C(i, j)$ , must have  $C(i-1, j-1)$ ,  $C(i, j-1)$ , and  $C(i-1, j)$  available.

Use  $(m+1) \times (n+1)$  matrix to store  $C$  values.

5. Backtrace for solution – either explicitly save indication of which of the three cases was used or recheck  $C$  values.

Takes  $\Theta(mn)$  time.

Example:

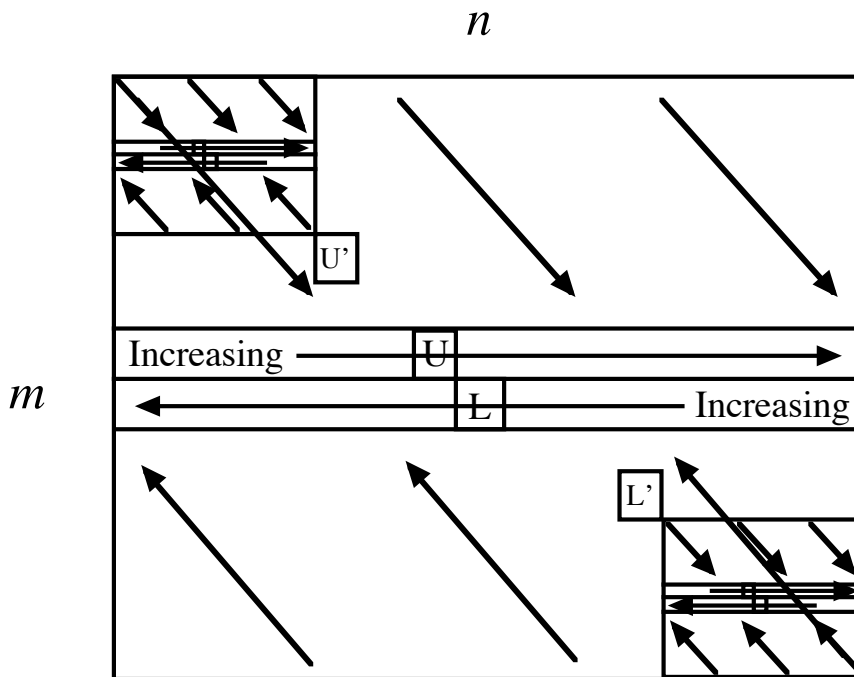
```

010101
001100
LCS is 0100, length==4
      0  0  1  1  0  0
      0  0  0  0  0  0  0
0  0  1  1  1  1  1  1
1  0  1  1  2  2  2  2
0  0  1  2  2  2  3  3
1  0  1  2  3  3  3  3
0  0  1  2  3  3  4  4
1  0  1  2  3  4  4  4

```

### Compact Version of Dynamic Programming

$\Omega(mn)$  space is *not* required.  $O(m+n)$  space is *attainable* using a recursive algorithm.



1. Use the usual NW to SE ( $\searrow$ ) cost computation to get increasing left-to-right cost row for LCS of the first  $m/2$  elements of the first sequence and the entire second sequence.
2. Perform symmetric (SE to NW,  $\swarrow$ ) cost computation for LCS of the last  $m/2$  elements of the first sequence and the entire second sequence.
3. Scan across the two seams for the diagonal pair (U,L) with the maximum sum. This sum is the length of the LCS.

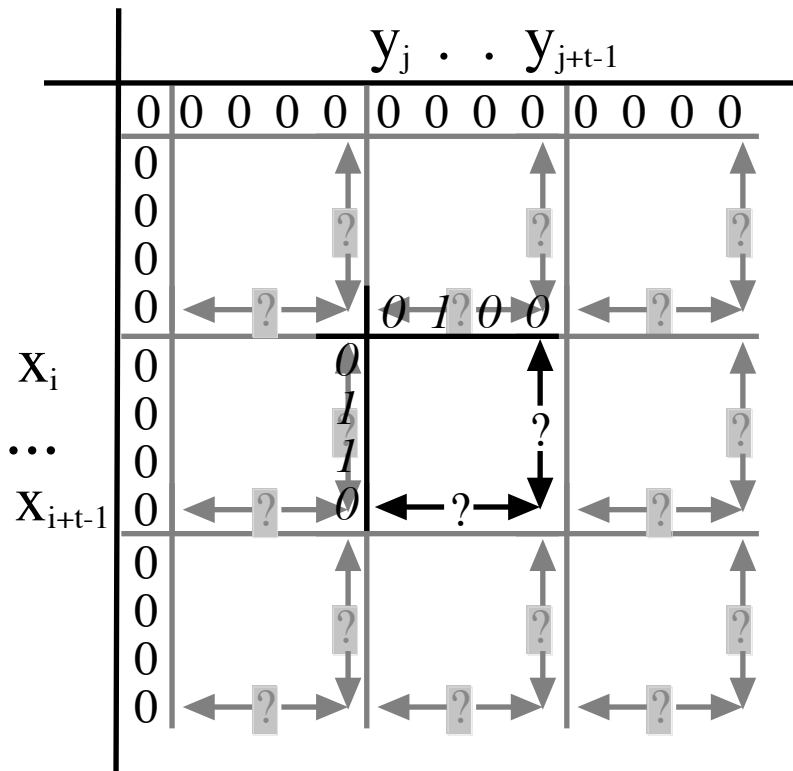


Four Russians' Method and LCS (aside, details in Gusfield)

When the LCS involves a small alphabet, the following properties allow Kronrod's technique to be applied:

1. Going across a row of the  $(m \times n)$  LCS matrix,  $C(i, j - 1) \leq C(i, j) \leq C(i, j - 1) + 1$ .
2. Going down a column of the LCS matrix,  $C(i - 1, j) \leq C(i, j) \leq C(i - 1, j) + 1$ .
3. Going down a diagonal of the LCS matrix,  $C(i - 1, j - 1) \leq C(i, j) \leq C(i - 1, j - 1) + 1$ .

A  $t \times t$  template is used to a) preprocess for all possible situations and b) tile the LCS matrix:



$0$  = value same as predecessor       $I$  = value is one more than predecessor

$x_i$  = position in first sequence       $y_j$  = position in second sequence

? = offset values supplied by preprocessing

Typically,  $t = \Theta(\log n)$  and runs in  $\Theta\left(\frac{n^2}{\log^2 n}\right)$ .

## LONGEST INCREASING SUBSEQUENCE

Equivalent to LCS:

Monotone: Sort input sequence A to obtain sequence A' in ascending order.

Compute LCS of A and A'.

Strict (LSIS): Remove duplicates from A' before computing LCS.

Dynamic Programming:

Concept: For each possible subsequence length, determine the smallest possible last element.

Algorithm: Process A left-to-right.

Maintain table of current smallest elements for the subsequence lengths.

*This table is non-decreasing.*

Each element is processed by a binary search and then linking to predecessor.

Monotone: Find last table element  $\leq$  current sequence element.

Strict: Similar, but current sequence element is ignored if already in table.

Monotone: 

0	1	2	3	4	5	6	7	8	9
3	2	1	1	2	3	3	2	1	2

  
(predecessors)

Strict: 

0	1	2	3	4	5	6	7	8	9
3	2	1	1	2	3	3	2	1	2

  
(predecessors)

## SPARSE LCS USING LSIS

Suppose an LCS for the following sequences is desired:

	0	1	2	3	4	5	6	7
sequence 1:	A	B	C	D	A	B	C	D
sequence 2:	A	A	B	B	C	C	D	D

1. Replace each occurrence of a symbol in sequence 1 by the positions of that symbol in sequence 2, but in *descending* order.

sequence 1:	A	B	C	D	A	B	C	D								
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
sequence 2 positions (predecessors)	1	<u>0</u>	3	2	5	4	7	6	<u>1</u>	0	3	<u>2</u>	5	<u>4</u>	7	<u>6</u>

2. Solve LSIS using dynamic programming.

Output gives positions in sequence 2 (0 1 2 4 6) for LCS elements: A A B C D

Takes time in  $O(r \log r)$  where  $r$  is the length of the constructed LSIS instance.

Binary search table may be replaced by van Emde Boas tree.