

4. Șiruri de caractere

- Marea majoritate a celor preocupați de activitatea de programare sunt familiarizați cu șirurile de caractere întrucât aproape toate limbajele de programare includ **șirul** sau **caracterul** printre elementele predefinite ale limbajului.
- În prima parte a capitolului se va preciza **tipul de date abstract șir**.
- Apoi vor fi abordate **modalitățile majore de implementare**, prin intermediul **tablourilor** respectiv al **pointerilor**.
- În ultima parte a capitolului vor fi precizate câteva din **tehnicele de căutare** în șiruri.

4.1. Tipul de date abstract șir

- Un șir este o colecție de caractere, spre exemplu "Structuri de date".
- În toate limbajele de programare în care sunt definite șiruri, acestea au la bază **tipul primitiv** `char`, care în afara literelor și cifrelor cuprinde și o serie de alte caractere.
- Se subliniază faptul că într-un șir, **ordinea caracterelor** contează. Astfel șirurile "CAL" și "LAC" deși conțin aceleași caractere sunt diferite.
- De asemenea, printr-un ușor **abuz** de notație se consideră că un *caracter* este interschimbabil cu *un șir* constând dintr-un singur caracter, deși strict vorbind ele sunt de tipuri diferite.
- Asemeni oricărui tip de date abstracte, definiția precisă a **TDA șir** necesită:
 - Descrierea **modelului** său **matematic**,
 - Precizarea **operatorilor** care acționează asupra elementelor tipului.
- Din punct de vedere **matematic**, elementele tipului de date abstract șir pot fi definite ca **secvențe finite de caractere** (c_1, c_2, \dots, c_n) unde c_i este de tip caracter, iar n precizează lungimea secvenței.
- Cazul în care n este egal cu zero, desemnează **șirul vid**.
- În continuare se prezintă un posibil set de operatori care acționează asupra TDA șir.
- Ca și în cazul altor structuri de date, există practic o libertate deplină în selectarea acestor operatori motiv pentru care setul prezentat are un caracter orientativ.

TDA Șir

Modelul matematic: secvență finită de caractere.

Notății: s, sub, u - șiruri;
 c - valoare de tip caracter;
 b - valoare booleană;
 $poz, lung$ - întregi pozitivi. [4.1.a]

Operatori:

CreazaSirVid(*s*) - procedură ce creează șirul vid *s*;
b:=SirVid(*s*) - funcție ce returnează true dacă șirul este vid;
b:=SirComplet(*s*) - funcție booleană ce returnează valoarea true dacă șirul este complet;
lung:=LungSir(*s*) - funcție care returnează numărul de caractere ale lui *s*;
poz:=PozitieSubsir(*sub,s*) - funcție care returnează poziția la care subșirul *sub* apare prima dată în *s*. Dacă *sub* nu e găsit în *s*, se returnează valoarea 0. Pozițiile caracterelor sunt numerotate de la stânga la dreapta începând cu 1;
ConcatSir(*u,s*) - procedură care concatenează la sfârșitul lui *u* atâtea caractere din *s*, până când **SirComplet**(*u*) devine **true**;
CopiazSubsir(*u,s,poz,lung*) - procedură care-l face pe *u* copia subșirului din *s* începând cu poziția *poz*, pe lungime *lung* sau până la sfârșitul lui *s*; dacă *poz>LungSir*(*s*) sau *poz<1*, *u* devine șirul vid;
StergeSir(*s,poz,lung*) - procedură care șterge din *s*, începând cu poziția *poz*, subșirul de *lung* caractere. Dacă *poz* este invalid (nu aparține șirului), *s* rămâne nemodificat;
InsereazaSir(*sub,s,poz*) - procedură care inserează în *s*, începând de la poziția *poz*, șirul *sub*;
c:=FurnizeazaCarSir(*s,poz*) - funcție care returnează caracterul din poziția *poz* a lui *s*. Pentru *poz* invalid, se returnează caracterul nul;
AdaugaCarSir(*s,c*) - procedură care adaugă caracterul *c* la sfârșitul șirului *s*;
StergeSubsir(*sub,s,poz*) - procedură care șterge prima apariție a subșirului *sub* în șirul *s* și returnează poziția *poz* de ștergere. Dacă *sub* nu este găsit, *s* rămâne nemodificat iar *poz* este poziționat pe 0;
StergeToateSubsir(*s,sub*) - șterge toate aparițiile lui *sub* în *s*.

- Operatorii definiți pentru un TDA-șir pot fi împărțiți în două categorii:
 - Operatori **primitivi** care reprezintă un set minimal de operații strict necesare în termenii cărora pot fi dezvoltate operatorii nonprimitivi.
 - Operatori **nonprimitivi** care pot fi dezvoltate din cei anteriori.
- Această divizare este oarecum **formală** deoarece, de obicei e mai ușor să definești un operator neprimitiv direct decât să-l definești în termenii primitivelor.
 - Spre **exemplu** operatorul **InsereazaSir** poate fi definit în termenii operatorilor **CreazaSirVid** și **AdaugaCar**.
 - Algoritmul este simplu: se construiește un șir de ieșire temporar (rezultat) căruia i se adaugă pe rând:
 - (1) caracterele șirului sursă până la punctul de inserție
 - (2) toate caracterele șirului de inserat (subșir)
 - (3) toate caracterele șirului sursă de după punctul de inserție.

- Șirul astfel construit înlocuiește șirul inițial (sursa) (secvența [4.1.b]).

```
/*Implementarea operatorului Insearea_Șir utilizând operatorii
Creaza_Sir_Vid, Adauga_Car_Sir și Furnizeaza_Car_Sir */
```

```
void insearea_sir(tipsir subsir, tipsir* sursa,tipindice p)
/*insearea subsirul în sursa între pozițiile p și p+1*/
{
    tipsir rezultat; tipindice i;

    if ((p<1) || (p>lung_sir(*sursa))) ;
        /*eroare (poziție ilegală în inserție)*/
    else
        {
            /*[4.1.b]*/
            creaza_sir_vid(rezultat);
            for( i= 1; i <= p-1; i++)
                adauga_car_sir(rezultat,
furnizeaza_car_sir(*sursa,i));
            for( i=1; i <= lung_sir(subsir); i++)
                adauga_car_sir(rezultat,
furnizeaza_car_sir(subsir,i));
            for( i=p; i <= lung_sir(*sursa); i++)
                adauga_car_sir(rezultat,
furnizeaza_car_sir(*sursa,i));
            creaza_sir_vid(*sursa);
            for( i=1; i <= lung_sir(rezultat); i++)
                adauga_car_sir(*sursa,
furnizeaza_car_sir(rezultat,i));
        }
}
```

4.2. Implementarea tipului de date abstract șir

- Sunt cunoscute două tehnici majore de implementare a șirurilor:
 - Implementarea bazată pe **tablouri**
 - Implementarea bazată pe **pointeri**

4.2.1. Implementarea șirurilor cu ajutorul tablourilor

- Cea mai simplă implementare a TDA-șir se bazează pe două elemente:
 - (1) Un **întreg** reprezentând lungimea șirului
 - (2) Un **tablou** de caractere care conține șirul propriu-zis.
- Un exemplu de astfel de implementare este următorul [4.2.1.a].

```
/* Exemplu de implementare a TDA Sir utilizând tablouri */
```

```
enum { lungime_max = 100};
typedef unsigned char tiplungime;

typedef unsigned char tipindice;

/*[4.2.1.a]*/
```

```
typedef struct tipsir {
    tiplungime lungime;
    char sir[lungime_max];
} tipsir;
tipsir s;
```

- Acest mod de implementare al șirurilor **nu** este unic.
- Se utilizează **tablouri** întrucât tablourile ca și șirurile sunt **structuri liniare**.
- Câmpul `lungime` este utilizat deoarece șirurile au lungimi diferite în schimb ce tablourile au lungimea fixă.
- Implementarea se poate dispersa de câmpul `lungime`, caz în care se poate utiliza un caracter convenit pe post de **santină de sfârșit (marker)**.
 - În această situație operatorul `LungimeȘir` va trebui să contorizeze într-o manieră liniară caracterele până la detectarea markerului.
 - Din acest motiv este preferabil ca lungimea să fie considerată un element **separat** al implementării.
- În contextul implementării șirurilor cu ajutorul tablourilor se definește noțiunea de **șir complet**
 - **Șirul complet** este șirul care are lungimea egală cu `lungime_maximă`, adică egală cu **dimensiunea** tabloului definit spre a-l implementa.
 - Șirurile implementate în acest mod **nu** pot depăși această lungime, motiv pentru care în acest context operează operatorul boolean `Șir_Complet`.
- În accepțiunea modelului anterior prezentat, în continuare se prezintă o implementare a operatorilor primitivi [4.2.1.b,c,d,e].

```
void creaza_sir_vid(tipsir* s) /*01*/
{
    s->lungime= 0; /*[4.2.1.b]*/
}
```

```
tiplungime lung_sir(tipsir* s) /*01*/
{
    tiplungime lung_sir_result;
    lung_sir_result= s->lungime; /*[4.2.1.c]*/
    return lung_sir_result;
}
```

```
char furnizeaza_car(tipsir s,tipindice poz) /*01*/
{
    char furnizeaza_car_result;
    if ((poz<1) || (poz>s.lungime)) /*[4.2.1.d]*/
    {
        /*eroare(poziție incorectă);*/
        furnizeaza_car_result= '/0'; /*caracterul vid*/
    }
    else
        furnizeaza_car_result= s.sir[poz-1];
    return furnizeaza_car_result;
}
```

```

-----
void adauga_car_sir(tipsir* s, char c)           /*O1*/
{
    if (s->lungime==lungime_max) ;
        /*eroare(se depașește lungimea maximă a șirului)*/
    else
    {
        s->lungime= s->lungime+1;           /*[4.2.1.e]*/
        s->sir[s->lungime-1]= c;
    }
}
-----

```

- Se observă că toate aceste operații rulează în $O(1)$ unități de timp indiferent de lungimea șirului.
- Procedurile CopiazăSubȘir, ConcatȘir, ȘtergeȘir și InsereazaȘir se execută într-un interval de timp liniar $O(n)$, unde n este după caz lungimea subșirului sau a șirului de caractere.
 - Spre **exemplu** procedura CopiazăSubșir ($u, s, poz, lung$) returnează în u subșirul având $lung$ caractere începând cu poziția poz .
 - Accesul la elementele subșirului se realizează direct ($s.șir[poz], s.șir[poz+1], \dots, s.șir[poz+lung-1]$), astfel încât consumul de timp al execuției este dominat de mutarea caracterelor [4.2.1.f].

```

-----
/*Implementarea operatorului CopiazăSubșir*/   /*O(n)*/

void copiazasubsir(tipsir* u, tipsir* s, tiplungime poz,
tiplungime lung)
{
    tiplungime indexsursa,indexcopiere;

    if ((poz<1) || (poz>s->lungime))
        u->lungime= 0;
    else
    {
        /*[4.2.1.f]*/
        indexsursa= poz-1;
        indexcopiere= 0;
        while ((indexsursa<s->lungime) &&
            (indexcopiere<u->lungime)) {
            indexsursa= indexsursa+1;
            indexcopiere= indexcopiere+1;
            u->sir[indexcopiere-1]= s->sir[indexsursa-1];
        } /*WHILE*/
        u->lungime= indexcopiere;
    }
}
-----

```

- Se observă faptul că în interiorul buclei **WHILE** există 3 atribuiri, care pentru o lungime $lung$ a sub șirului determină $3*lung + 3$ atribuiri.
- Se observă de asemenea că procedura CopiazăSubșir este liniară în raport cu lungimea $lung$ a subșirului.
- Un alt exemplu îl reprezintă implementarea funcției PozițieSubșir [4.2.1.g].

```

-----
/*Implementarea operatorului PozițieSubșir*/

tiplungime pozitiewsubsir(tipsir* sub,tipsir* s)
{
    tipindex mark, /*index pentru punctul de start al unei
comparatii*/
    indexsub,      /*index subsir*/
    indexsir;      /*index sir*/
    tiplungime poz; /*pozitia lui sub în s*/
    boolean stop;  /*devine adevarat când elementele
                    corespozatoare din s si sub nu sunt egale*/

    tiplungime pozitiewsubsir_result;
    indexsir= 1;
    poz= 0;
    do {
        indexsub= 1;
        mark= indexsir;
        stop= false;
        while ((! stop) && (indexsir<=s->lungime) &&
                (indexsub<=sub->lungime))
            if (s->sir[indexsir-1]==sub->sir[indexsub-1])
                {
                    indexsir= indexsir+1;
                    indexsub= indexsub+1;          /*[4.2.1.g]*/
                }
            else
                stop= true; /*WHILE*/
        if (indexsub>sub->lungime)
            poz= mark; /*potrivire*/
        else
            indexsir= mark+1;
    } while (!((poz>0) ||
                (indexsir+sub->lungime-1>s->lungime)));
    pozitiewsubsir_result= poz;
    return pozitiewsubsir_result;
}
-----

```

- Complexitatea funcției PozițieSubșir(sub,s) este $O(\text{lungime} * s.\text{lungime})$ unde lungime este lungimea modelului iar s.lungime este lungimea șirului în care se face căutarea.
- Există însă și alte metode mult mai **performante** de căutare care fac obiectul unui paragraf ulterior.

4.3. Tehnici de căutare în șiruri

- Una din operațiile cele mai frecvente care se execută asupra șirurilor este **căutarea**.
- Întrucât **performanța** acestei operații are o importanță covârșitoare asupra mării majorități a prelucrărilor care se realizează într-un sistem de calcul, studiul **tehnichilor de căutare performante** reprezintă o **preocupare permanentă** a cercetătorilor în domeniul programării.

- În cadrul acestui paragraf vor fi trecute în revistă câteva dintre cele mai cunoscute tehnici de căutare în șiruri de caractere.

4.3.1. Căutarea de șiruri directă

- O manieră frecvent întâlnită de căutare este așa numită **căutare de șiruri directă** ("string search").
- **Formularea problemei:**
 - Se dă un tablou s de n elemente și un tablou p de m elemente, unde $0 < m < n$ declarate astfel [4.3.2.a]:

```
-----
typedef unsigned char boolean;
char s[n];
char p[m];                                     /* [4.3.2.a] */
-----
```

- Căutarea se șiruri are drept scop stabilirea **primei apariții** a lui p în s .
- De regulă, s poate fi privit ca un **text**, iar p ca un cuvânt **model** ("pattern") a cărui **primă apariție** se caută în textul s .
- Aceasta este o operație **fundamentală** în toate sistemele de prelucrare a textelor și în acest sens este de mare interes găsirea unor algoritmi cât mai eficienți.
- Cea mai **simplă** metodă de căutare este așa numita **căutare de șiruri directă**.
- **Rezultatul** unei astfel de căutări este un indice i care precizează apariția unei **coincidențe** între model și șir.
- Acest lucru este formalizat de predicatul $P(i, j)$ [4.3.2.b].

$$P(i, j) \leq A_k : 0 \leq k < j : s_{i+k} = p_k \quad [4.3.2.b]$$

- Predicatul $P(i, j)$ precizează faptul că există o **coincidență** între:
 - Șirul s (începând cu indicele i)
 - Șirul p (începând cu indicele 0) pe o lungime de j caractere (fig. 4.3.2.a).

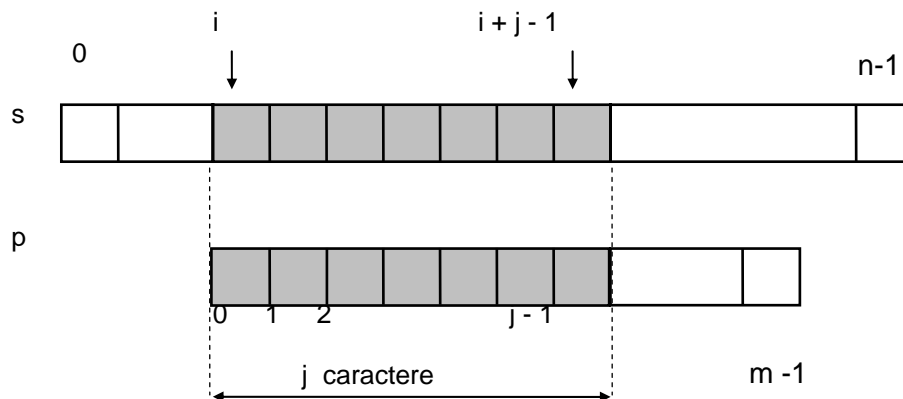


Fig.4.3.2.a. Reprezentarea grafică a predicatului $P(i, j)$

- Este evident că indexul i care rezultă din **căutarea directă de șiruri** trebuie să satisfacă predicatul $P(i, m)$.
- Această condiție **nu** este însă suficientă.
- Deoarece căutarea trebuie să furnizeze **prima** apariție a modelului, $P(k, m)$ trebuie să fie fals pentru toți $k < i$.
- Se notează această condiție cu $Q(i)$ [4.3.2.c].

$$Q(i) = \bigwedge_{k: 0 \leq k < i} \sim P(k, m) \quad [4.3.2.c]$$

- Punerea problemei sugerează formularea căutării directe de șiruri ca și o iterație de comparații redactată în **termenii** predicatelor Q respectiv P .
- Astfel implementarea lui $Q(i)$ conduce la secvența [4.3.2.d]:

```
{Implementarea predicatului Q(i)}
i:= -1;
REPEAT
    i:= i+1;
    gasit:= P(i,m)
UNTIL gasit OR (i=n-m);
```

- Calculul lui P rezultă în continuare ca și o iterație de comparații de caractere individuale.
- Rafinarea secvenței anterioare conduce de fapt la **implementarea căutării directe de șiruri** ca o repetiție într-o altă repetiție [4.3.2.e].

```
/*{Implementarea căutării directe de șiruri*/

boolean cautaredirecta(int* poz)
{
    int i,j;    /*i parcurge caracterele din sir,
                j parcurge caracterele din model*/
    boolean cautaredirecta_result;
    i= 1;
    do {
        i= i+1; j= 0;    /*Q(i)*/
        while ((j<m) && (s[i+j]==p[j]))
            j= j+1;    /*P(i,m)*/
    } while (!(j==m || i==n-m));
    *poz=i;
    cautaredirecta_result= j==m;
    return cautaredirecta_result;
}
```

- Termenul $j = m$ din condiția de terminare, corespunde lui **găsit** deoarece el implică $P(i, m)$.
- Termenul $i = n - m$ implică $Q(n - m)$, deci **non** existența vreunei coincidențe în cadrul șirului.
- **Analiza căutării de șiruri directe.**

- Algoritmul lucrează destul de **eficient** dacă se presupune că nepotrivirea în procesul de căutare apare cel mult după câteva comparații în cadrul buclei interioare.
- Astfel pentru un set de 128 de caractere se poate presupune că nepotrivirea apare după inspectarea a 1 sau 2 caractere.
- Cu toate acestea în cazul cel mai nefavorabil, degradarea performanței este îngrijorătoare.
 - Astfel dacă spre exemplu șirul s este format din $n-1$ caractere 'A' urmate de un singur 'B',
 - Iar modelul constă din $m-1$ caractere 'A' urmate de un 'B',
 - În acest caz este necesar un număr de comparații de ordinul $n * m$ pentru a găsi coincidența la sfârșitul șirului.
 - Din fericire există metode care îmbunătățesc radical comportarea algoritmului în această situație.
 - Tehnicile de căutare care sunt prezentate în continuare materializează acest deziderat.

4.3.2. Căutarea de șiruri Boyer-Moore

- Metoda de căutare inventată în 1975 de R.S. Boyer și J.S. Moore este una dintre cele mai ingenioase și în același timp una dintre cele mai performante.
- Căutarea BM, după cum mai este numită, este bazată pe ideea **neobișnuită** de a începe compararea caracterelor de la **sfârșitul** modelului și **nu** de la început.
- Modelul este **precompilat** în prealabil într-un tablou d .
 - Astfel, pentru fiecare caracter x al setului de caractere, se notează cu d_x distanța dintre cea mai din dreapta apariție a lui x în cadrul modelului și sfârșitul modelului (fig.4.3.4.a.).

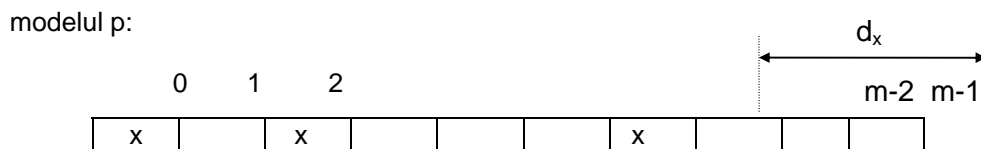


Fig. 4.3.4.a. Determinarea valorii d_x corespunzătoare caracterului x al modelului

- Valoarea d_x se trece în tabloul d în poziția corespunzătoare caracterului x .
- Pentru caracterele care **nu** sunt în model respectiv pentru ultimul caracter al modelului, d_x se face egal cu **lungimea totală** a modelului.
- În continuare, se presupune că în procesul de **comparare** al șirului cu modelul a apărut o **nepotrivire** între caracterele corespunzătoare.
 - În această situație modelul poate fi imediat **deplasat** spre dreapta cu $d[p_{m-1}]$ poziții, valoare care este de regulă mai mare ca 1.

- Se precizează faptul că p_{m-1} este caracterul din șirul baleat s , corespunzător **ultimului** caracter al modelului la momentul considerat, **indiferent** de locul în care s-a constatat nepotrivirea (figura 4.3.4.b.).

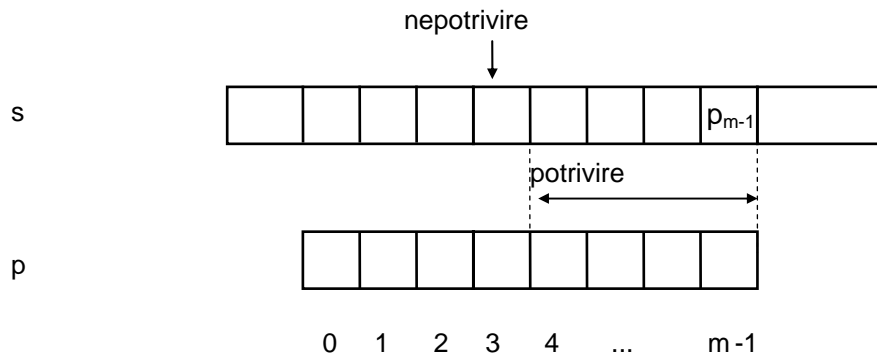
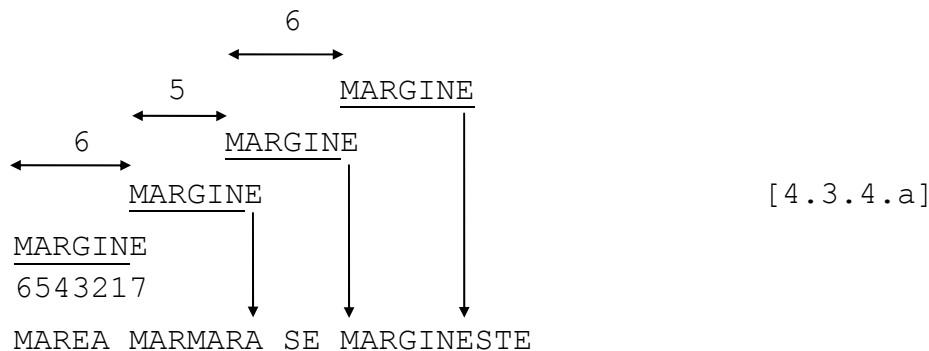


Fig. 4.3.4.b. Comparare șir - model pentru determinarea lui p_{m-1}

- Dacă caracterul p_{m-1} **nu** apare în model, deplasarea este mai mare și anume cu întreaga lungime a modelului.
- Exemplul din secvența [4.3.4.a.] evidențiază acest proces.



- Implementarea căutării Boyer-Moore apare în secvența [4.3.4.c.]

```

/*căutare Boyer-Moore */
i= m; j= m;
while ((j>0) && (i<n))
{
    /*Q(i-m)*/
    j= m; k= i;
    while ((j>0) && (s[k-1]==p[j-1]))
    {
        /*P(k-j, j) & (k-j=i-m) */
        k= k-1;
        j= j-1;
    }
    i= i+d[s[i-1]];
}
/*cautare Boyer-Moore*/

```

- Indicii implicați satisfac următoarele relații: $0 < j < m$ și $0 < i, k < n$.

- Astfel **terminarea** algoritmului cu $j = 0$ implică o **potrivire** începând de la poziția k spre dreapta, de m caractere, unde $k = i - m$.
- **Terminarea** algoritmului cu $j > 0$ implică $i = n$ element care indică **absența** potrivirii.
- Programul următor [4.3.4.d] implementează strategia Boyer-Moore într-un context mai general.

```

/* Căutarea Boyer-Moore -varianta C */

typedef unsigned char boolean;
#define true          (1)
#define false         (0)

enum { mmax = 100 /*lungime maxima model*/,
       nmax = 200} /*lungime maxima sir sursa*/; /*[4.3.4.c]*/
int m/*lungime model*/,n/*lungime sir*/;
char p[mmax];          /*model*/
char s[nmax];          /*sir*/
int d[256];            /*tabela de deplasari*/
boolean cautarebm(int* poz)
{
    int i,j,k;

    /*citire sir;    n este lungimea curenta a sirului
    /*citire model;  m este lungimea curenta a modelului*/
    boolean cautarebm_result;
    for( i=0; i <= 255; i++) d['i']= m; /*compilare model*/
    for( j=0; j <= m-2; j++) d[p[j]]= m-j-1;
    i=m; j=m; /*cautare model*/
    while ((j>0) && (i<=n))
    {
        j= m;
        k= i;
        while ((j>0) && (s[k-1]==p[j-1]))
        {
            k= k-1;
            j= j-1;
        }
        if (j>0)
            i=i+d[s[i-1]];
    }
    *poz=i-m; /*poz=k*/
    cautarebm_result= j==0;
    return cautarebm_result;
}

```

• Analiza căutării Boyer-Moore

- Autorii căutării BM, au demonstrat proprietatea remarcabilă că în **toate** cazurile, cu excepția unora special construite, **numărul de comparații** este substanțial mai **redus** decât n .
- În cazul cel mai **favorabil** când ultimul caracter al modelului nimereste întotdeauna în șir peste un caracter diferit de cele ale modelului, **numărul de comparații** este n/m .

4.4. Rezumat

- Un **șir** este o colecție de caractere. În toate limbajele de programare în care sunt definite șiruri, acestea au la bază tipul primitiv `char`.
- Asemeni oricărui tip de date abstracte, **definirea TDA șir** necesită pe de o parte descrierea **modelului său matematic**, iar pe de altă parte precizarea **operatorilor** care acționează asupra elementelor tipului.
- Cea mai cunoscută modalitate de implementare a TDA șir o constituie **tablourile liniare**.
- Dintre operațiile cele mai frecvente care se execută asupra șirurilor de caractere este **căutarea**.
- Cea mai frecventă metodă de căutare este **căutarea de șiruri directă ("string search")**.
- O metodă de căutare cu performanțe mult mai bune este **căutarea de șiruri Boyer-Moore**.

4.5. Exerciții

- 1) Se cere să se precizeze *modelul matematic* și *setul de operatori* pentru *TDA șir de caractere*.
- 2) Se cere să se realizeze o implementare bazată pe tablouri a următorilor operatori: *CreazăȘirVid*, *LungȘir*, *InsereazăȘir*, *ConcatȘir*, *PozițieSubșir*.
- 3) Se cere să se redacteze o funcție care implementează *căutarea directă de șiruri*.
- 4) Care este *principiul căutării Boyer-Moore*? Cum se realizează *precompilarea modelului*?
- 5) Se cere să se redacteze o funcție care implementează *căutarea de șiruri Boyer-Moore*.
- 6) Se cere să se realizeze un program C care:
 - citește un șir de caractere `s` de la tastatură
 - citește un șir de caractere model `p`
 - caută în șirul `s` pe `p` utilizând căutarea directă de șiruri
 - caută în șirul `s` pe `p` utilizând căutarea de șiruri BM

Obs. Se vor utiliza funcțiile dezvoltate în exercițiile anterioare.