

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**.
 - În continuare vor fi abordate **modalitățile majore de implementare a șirurilor**, prin intermediul **tablourilor** respectiv al **pointerilor**.
 - În ultima parte a capitolului vor fi precizate câteva din **tehnicile 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, definirea precisă a **TDA Șir** necesită:
 - (1) Descrierea **modelului** său **matematic**.
 - (2) 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. [4.1.a]

Notatii: *s, sub, u* - siruri;
c - valoare de tip caracter;
b - valoare booleană;
poz, lung - întregi pozitivi.

Operatori:

CreazaSirVid(*sir s*) - procedură ce creează șirul vid *s*;

b=**SirVid**(*sir s*) - funcție ce returnează true dacă șirul *s* este vid;

b=**SirComplet**(*sir s*) - funcție booleană ce returnează valoarea true dacă șirul este complet;

lung=**LungSir**(*sir s*) - funcție care returnează numărul de caractere ale lui *s*;

poz=**PozitieSubsir**(*sir 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(*sir u, s*) - procedură care concatenează la sfârșitul lui *u* atâtea caractere din *s*, până când **SirComplet**(*u*) devine **true**;

CopiazaSubsir(*sir u, s, int 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(*sir s, int 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(*sir sub, s, int poz*) - procedură care inserează în *s*, începând de la poziția *poz*, șirul *sub*; eventual trunchiere la depășirea lungimii maxime;

c=FurnizeazaCarSir(sir s, int poz) - funcție care returnează caracterul din poziția *poz* a lui *s*. Pentru *poz* invalid, se returnează caracterul nul;

AdaugaCarSir(sir s, char c) - procedură care adaugă caracterul *c* la sfârșitul șirului *s*;

StergeSubsir(sir sub,s, int 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(sir 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 dezvoltati operatorii nonprimitivi.
 - Operatori **nonprimitivi** care pot fi dezvoltati 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 unor operatori primitivi.
 - Spre **exemplu** operatorul **InseereazăȘir** poate fi definit în termenii operatorilor **CreazăȘirVid**, **AdaugăCar** și **FurnizeazaCar**.
 - 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 Inseerează_Șir utilizând
operatorii Crează_Sir_Vid, Adaugă_Car_Sir și
Furnizeaza_Car_Sir */
                                                                    /*[4.1.b]*/
void inseereaza_sir(tip_sir subsir, tip_sir* sursa, int p)
/*inseerează subșir în sursă între pozițiile p și p+1*/
{
    tip_sir rezultat; tip_indice i;

    if ((p<1) || (p>lung_sir(*sursa))) ;
    /*eroare (poziție ilegală în inserție)*/
```

```

else
{
    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));
}
}{insereaza_sir}
-----

```

4.2. Implementarea tipului de date abstract şir

- Sunt cunoscute **două tehnici majore** de implementare a şirurilor:
 - (1) Implementarea bazată pe **tablouri**
 - (2) Implementarea bazată pe **pointeri**.

4.2.1. Implementarea şirurilor cu ajutorul tablourilor

- Cea mai simplă implementare a **TDA-şir** se bazează pe două elemente:
 - (1) Un **întreg** reprezentând lungimea şirului.
 - (2) Un **tablou** de caractere care conţine şirul propriu-zis.
- În tablou caracterele pot fi păstrate ca atare sau într-o formă **împachetată**.
- Exemple pentru o astfel de implementare în Pascal respectiv C apar în secvenţa [4.2.1.a].

```

/* Exemplu de implementare a TDA Sir utilizând tablouri -
structuri de date - varianta C*/

```

```

enum {Lungime_Max=100};
/*[4.2.1.a]*/

typedef struct tip_sir {
    int lungime;          /*lungime curentă şir*/
    char sir[Lungime_Max]; /*tablou suport şir*/
} tip_sir;

```

```
tip_sir s;                                /*instantiere şir*/
```

- Acest mod de implementare al şirurilor **nu** este unic.
- Se utilizează **tablouri** întrucât tablourile ca şi şirurile sunt **structuri liniare**.
 - Câmpul `lungime` este utilizat deoarece şirurile 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 **santinelă de sfârşit (marker)**.
 - În această situaţie operatorul `LungimeŞir` va trebui să **contorizeze** într-o manieră liniară caracterele până la detectarea markerului.
 - Din acest motiv este preferabil ca **lungimea şirului** să fie considerată un element **separat** al implementării.
- În contextul implementării şirurilor cu ajutorul tablourilor se defineşte noţiunea de **şir complet**.
 - **Şirul complet** este şirul care are lungimea egală cu `Lungime_Max`, adică egală cu **dimensiunea** tabloului definit spre a-l implementa.
 - Şirurile implementate în acest mod **nu** pot depăşi această lungime, motiv pentru care în acest context operează operatorul boolean **`ŞirComplet`**.
- În accepţiunea modelului anterior prezentat, în continuare se prezintă o implementare a operatorilor primitivi [4.2.1.b,c,d,e] varianta C.

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

```
void creaza_sir_vid(tip_sir* s)                /*[4.2.1.b]*/  
{  
    s->lungime= 0;                            /*performanţa O(1)*/  
} /*creaza_sir_vid*/
```

```
-----  
int lung_sir(tip_sir* s)                     /*[4.2.1.c]*/  
{  
    int lung_sir_result;                      /*performanţa O(1)*/  
    lung_sir_result=s->lungime;  
    return lung_sir_result;  
} /*lung_sir*/
```

```
-----  
char furnizeaza_car(tip_sir s,tip_indice poz) /*[4.2.1.d]*/  
{  
    char furnizeaza_car_result;  
    if ((poz<1) || (poz>s.lungime))  
    {  
        /*eroare(pozitie incorectă);*/  
        furnizeaza_car_result='/0'; /*caracterul vid*/  
    }  
    Else                                /*performanţa O(1)*/  
        furnizeaza_car_result=s.sir[poz-1];
```

```

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

```

- Se observă că toate aceste operații rulează în $O(1)$ unități de timp indiferent de lungimea șirului.
- Procedurile **CopiazăSubșir**, **Concatșir**, **Ștergeșir** și **Inseereazașir** 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 din s având lung caractere începând cu poziția poz.
 - Accesul la elementele subșirului se realizează direct (s.șir[poz], s.șir[poz+1],...,s.șir[poz+lung-1]), astfel încât consumul de timp al execuției este dominat de mutarea caracterelor [4.2.1.f].

```

/*Implementarea operatorului CopiazăSubșir*/

void copiaza_subsir(tip_sir* u, tip_sir* s, int poz, int
                    lung)
/*copiază un subșir din s, de la poziția poz, de lungime
lung în șirul u*/
{
    /*[4.2.1.f]*/
    int index_sursa,index_copiere;

    if ((poz<1) || (poz>s->lungime))
        u->lungime=0;
    else
        /*performanța O(n)*/
        {
            index_sursa=poz-1;
            index_copiere=0;
            while ((index_sursa<s->lungime) &&
                    (index_copiere<lung))
            {
                index_sursa=index_sursa+1;
                index_copiere=index_copiere+1;
                u->sir[index_copiere-1]=s->sir[index_sursa-1];
            } /*while*/
            u->lungime=index_copiere;
        } /*else*/
}

```

```
}/*copiaza_subsir*/
```

- Se observă faptul că în interiorul buclei **WHILE** există 3 atribuiri, care pentru o lungime `lung` a sub şirului determină $3 \cdot \text{lung} + 3$ atribuiri.
- Se observă de asemenea că procedura `CopiazăSubşir` este liniară în raport cu lungimea `lung` a subşirului.
- Un alt exemplu îl reprezintă implementarea funcţiei `PoziţieSubşir` [4.2.1.g].
- Procedura de lucru a operatorului ***pozitie_subsir*** este următoarea:
 - Se compară şirul original `s` cu şirul model căutat sub poziţie cu poziţie, până la determinarea primei neconcordanţe. Simultan se memorează în `mark` punctul de start în şirul `s` al comparaţiei curente. Dacă s-a găsit coincidenţa între `s` şi `sub` se returnează poziţia lui `sub` în `s` (începutul coincidenţei, adică valoarea lui `mark`). Dacă nu s-a găsit o coincidenţă, `sub` avansează în `s` cu o poziţie (`mark=mark+1`) şi procesul se reia până la găsirea unei coincidenţe, sau până când în şirul original `s` nu mai există suficiente caractere pentru lungimea modelului `sub`. Căutarea nereuşită returnează valoarea `-1`.

```
/*Implementarea operatorului PoziţieSubşir*/
```

```
/*[4.2.1.g]*/
```

```
int pozitie_subsir(tip_sir* sub,tip_sir* s)  
/*determină poziţia lui sub în cadrul şirului s (prima  
apariţie) şi returnează poziţia găsită sau -1 în caz de  
nereuşită*/
```

```
{  
    int mark;          /*index pentru punctul de start al unei  
                        comparaţii*/  
    int index_sub; /*index subşir*/  
    int index_sir; /*index şir*/  
    int poz;          /*poziţia lui sub în s*/  
    boolean stop;      /*devine adevărat când elementele  
                        corespunzătoare din s şi sub nu sunt egale*/
```

```
    index_sir=0;  
    poz=-1;           /*inițializare valoare poz pe -1*/
```

```
    do {  
        index_sub=0;  
        mark= index_sir;  
        stop=false;  
        while((!stop)&&(index_sir<s->lungime)&&  
              (index_sub<sub->lungime))  
            if(s->sir[index_sir]==sub->sir[index_sub])  
            {  
                index_sir=index_sir+1;  
                index_sub=index_sub+1;  
            }  
        else  
            stop=true; /*while*/  
    }  
    if(index_sub>sub->lungime)
```

```

        poz=mark;  /*potrivire*/
    else
        index_sir=mark+1; /*nepotrivire => avans model*/
    } while (! (poz>0) ||
        (index_sir+sub->lungime-1>s->lungime)); /*do*/
    return poz;
} /*pozitie_subsir*/

```

- Complexitatea funcției **pozitie_subsir**(sub,s) este $O(sub.lungime * s.lungime)$ unde `sub.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.
- Unul dintre marile **avantaje** ale utilizării datelor abstracte este următorul:
 - În situația în care se găsește un algoritm mai **performant** pentru o anumită operație, se poate foarte simplu înlocui o **implementare** cu alta fără a afecta restul programului în condițiile păstrării nealterate a prototipului operatorului.

4.2.2. Tabele de șiruri

- Utilizarea tabloului în implementarea șirurilor are avantajul că operatorii sunt ușor de redactat, simplu de înțeles și suficient de rapizi.
- Cu toate că acest mod de implementare este economic din punct de vedere al **timpului** de execuție, el este ineficient din punctul de vedere al **spațiului** utilizat.
 - Lungimea **maximă** a tabloului trebuie aleasă la dimensiunea **celui mai lung șir** preconizat a fi utilizat deși în cea mai mare parte a cazurilor se utilizează șiruri mult mai **scurte** .
- O modalitate de a economisi spațiul în dauna vitezei de lucru specifică implementării cu ajutorul tablourilor, este aceea de a utiliza următoarele structuri de date:
 - (1) Un **tablou** suficient de lung numit " **heap** " (grămadă) pentru a memora toate șirurile.
 - (2) O **tabelă** auxiliară de șiruri care păstrează evidența șirurilor în heap.
- Aceasta tehnică este utilizată în unele interprete BASIC.
- În acest mod de implementare a șirurilor, un șir este reprezentat printr-o **intrare** în **tabela de șiruri** care conține două valori (fig.4.2.2.a):
 - **Lungimea** șirului .
 - Un **indicator** în heap care precizează **începutul** șirului.
- Un exemplu de implementare în limbajul C apare în secvența [4.2.2.a].
- Se observă că această implementare presupune:
 - Variabilele globale `TabelaDeȘiruri`, `Heap`, `Disponibil` și `NumărDeȘiruri`.

- Structurarea în forma unui articol cu câmpurile lungime și indicator a elementelor tablei de șiruri.

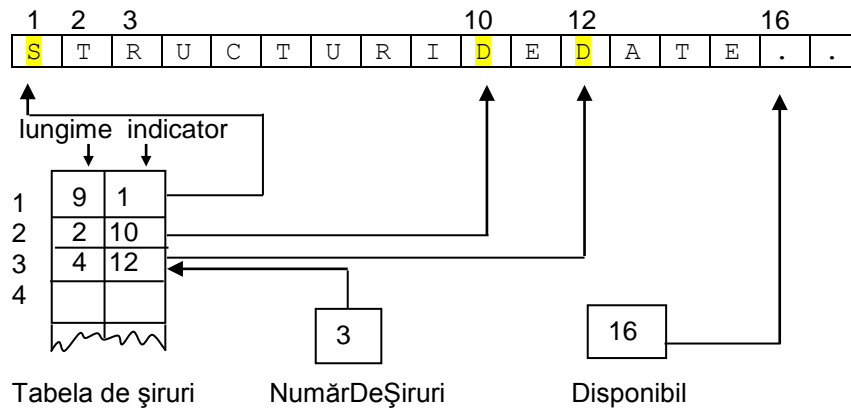


Fig.4.2.2.a. Modelul tablei de șiruri

/*Implementarea șirurilor utilizând tabele - structuri de date*/

/*[4.2.2.a]*/

int Lungime_Heap = **/*numărul maxim de caractere*/**;

int Lungime_Tabela = **/*numărul maxim de șiruri*/**;

typedef struct element_tabela

```
{
    int lungime;        /*lungimea șirului*/
    int indicator;      /*poziție în heap*/
} Element_Tabela;
```

typedef int Tip_Sir;

/*tabela de șiruri*/

element_tabela Tabela_De_Siruri[Lungime_Tabela];

char Heap[Lungime_Heap]; **/*zona de șiruri (heap)*/**

int Disponibil; **/*prima locație liberă în heap*/**

int Numar_De_Siruri; **/*număr curent de șiruri*/**

- În continuare se prezintă pentru exemplificare implementarea operatorului **AdaugăCarȘir**.

- Deoarece adăugarea unui caracter produce depășirea șirului în cauză în dauna șirului următor din spațiul de lucru, funcția **AdaugăCarȘir** trebuie:
 - [1] Să creeze o nouă copie a șirului în zona disponibilă a heap-ului.
 - [2] Să recopieze șirul în noua locație.
 - [3] Să adauge noul caracter.

- [4] Să reactualizeze tabela de șiruri și variabila Disponibil [4.2.2.b].

```
/*Implementarea operatorului AdaugăCarȘir - varianta pseudocod*/
```

```
PROCEDURE AdaugaCarSir(Tip_Sir s, char c) /*[4.2.2.b]*/
/*adaugă caracterul c la sfârșitul șirului s*/
  int i, lungimeVeche, indicatorVechi;

  daca ((s<1) OR (s>Numar_De_Siruri)) /*performanța O(n)*/
    *eroare(referință invalidă de șir);
  altfel
    [1] lungimeVeche=Tabela_De_Siruri[s].lungime;
        indicatorVechi=Tabela_De_Siruri[s].indicator;

    [2] pentru i=0 la lungimeVeche-1
        Heap[Disponibil+i]=Heap[indicatorVechi+i];

    [3] Heap[Disponibil+lungimeVeche]=c; /*se adaugă c*/

    [4] Tabela_De_Siruri[s].indicator=Disponibil;
        Tabela_De_Siruri[s].lungime=lungimeVeche+1;
        Disponibil=Disponibil+lungimeVeche+1;
        □ /*altfel*/
/*AdaugaCarSir*/
```

- Se observă că în această implementare **AdaugăCarȘir** necesită $O(n)$ unități de timp pentru un șir de n caractere față de implementarea bazată pe tablouri care necesită un timp constant ($O(1)$).
 - Aceasta este tributul plătit facilității de a putea lucra cu șiruri de **lungime variabilă**.
 - De asemenea în urma creării noii copii a șirului supus operației de adăugare, vechea instanță a șirului rămâne în heap ocupând memorie în mod inutil.
- O altă problemă potențială a acestei implementări este următoarea:
 - Se consideră o procedură care realizează **copierea** unui șir **sursă** într-un șir **destinație**.
 - În implementarea bazată pe **tabelă de șiruri** acest lucru se poate realiza simplu făcând ca șirului destinație să-i corespundă **același indicator** în heap ca și șirului sursă.
 - Acest fenomen este cunoscut în programare sub denumirea de **supraîncărcare** ("**aliasing**"), adică două variabile diferite se referă la aceeași locație de memorie.
 - Dacă se cere însă **ștergerea** șirului sursă sau destinație se pot pierde ambele șiruri.
 - Această problemă poate fi ocolită printr-o implementare adecvată.

- Pentru rezolvarea **problemei supraîncărcării**, se poate proceda ca și în cazul operatorului **AdaugăCarȘir**, adică:
 - (1) Pentru șirul destinație se creează o **instanță** începând cu prima locație disponibilă în heap.
 - (2) Se copiază sursa în noua locație.
 - (3) Se introduce referința care indică această locație, în tabela de șiruri în intrarea corespunzătoare șirului nou creat.
- De fapt, în anumite variante ale limbajului BASIC, această manieră de rezolvare a copierii stă la baza implementării instrucțiunii LET.

4.2.3. Implementarea șirurilor cu ajutorul pointerilor

- Esența implementării **șirurilor** cu ajutorul pointerilor rezidă în reprezentarea acestora printr-o colecție de **celule înlănțuite**.
- Fiecare **celulă** conține [4.2.3.a]:
 - Un **caracter** (sau un grup de caractere într-o implementare mai elaborată).
 - Un **pointer** la celula următoare.

 {Implementarea șirurilor cu ajutorul pointerilor - structuri de date - varianta PASCAL}

```
TYPE PointerCelula = ^Celula;                                {[4.2.3.a]}
```

```

Celula = RECORD
    ch: char;
    urm: PointerCelula
END;

TipSir = RECORD
    lungime: integer;
    cap,coada: PointerCelula
END;
```

 /*Implementarea șirurilor cu ajutorul pointerilor - structuri de date - varianta C*/

```
typedef struct celula                                        /*[4.2.3.a]*/
{
    char ch;
    struct celula* urm;
} celula;
```

```
typedef struct Tip_Sir
{
    int lungime;
```

```

    celula* cap,coada;
} Tip_Sir;

```

- Spre exemplu în fig.4.2.3.a apare reprezentarea şirului 'DATE' în această implementare,

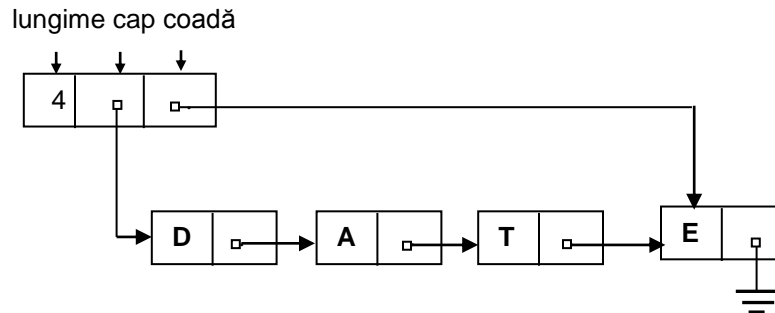


Fig. 4.2.3.a. Implementarea şirurilor cu ajutorul pointerilor

- În secvenţele [4.2.3.b, c, d, e] apare implementarea unor operatori primitivi specifici.

{Implementarea şirurilor cu ajutorul pointerilor}

```

PROCEDURE CreazaSirVid(VAR s: TipSir);           {[4.2.3.b]}
BEGIN
    s.lungime:=0;
    s.cap:=nil;                                   {performanţa O(1)}
    s.coada:=nil
END; {CreazaSirVid}

```

```

FUNCTION LungSir(s: TipSir): integer;           {[4.2.3.c]}
BEGIN
    LungSir:=s.lungime                           {performanţa O(1)}
END; {LungSir}

```

```

FUNCTION FurnizeazaCarSir(s:TipSir; poz:integer): char;
VAR contor: integer;
    p: PointerCelula;                            {[4.2.3.d]}
BEGIN
    IF (poz<1)AND(poz>s.lungime) THEN
        BEGIN
            *eroare(index eronat)
            FurnizeazaCarSir:=chr(0) {caracterul nul}
        END
    ELSE                                         {performanţa O(n)}
        BEGIN
            p:=s.cap;
            {căutare secvenţială în şir}
            FOR contor:=1 TO poz-1 DO p:=p^.urm;
            FurnizeazaCarSir:=p^.ch
        END
    END

```

```

END; {FurnizeazaCarSir}
-----
PROCEDURE AdaugaCarSir(var s:TipSir; c:char); {[4.2.3.e]}
BEGIN
  IF s.lungime=0 THEN {primul caracter al șirului}
    BEGIN
      new(s.cap);
      s.cap^.ch:=c;
      s.cap^.urm:=nil;
      s.coada:=s.cap
    END
  ELSE {performanța O(1)}
    BEGIN
      new(s.coada^.urm);
      s.coada^.urm^.ch:=c;
      s.coada^.urm^.urm:=nil;
      s.coada:=s.coada^.urm
    END; {ELSE}
  s.lungime:=s.lungime+1
END; {AdaugaCarSir}
-----

```

- În ceea ce privește implementarea operației de **copiere**, apar aceleași probleme legate de **supraîncărcare** care pot fi rezolvate într-o manieră asemănătoare prin generarea unui **șir nou destinație** identic cu sursa.
- Un astfel de operator **Copiază**(destinație, sursa) poate fi utilizat cu succes în implementarea operației **ConcatȘir** [4.2.3.f].

{Implementarea operatorului ConcatȘir utilizând operatorul Copiază}

```

PROCEDURE ConcatSir(VAR u: TipSir; s: TipSir); {[4.2.3.f]}
VAR temp: TipSir;
BEGIN
  IF u.lungime=0 THEN
    Copiază(u, s)
  ELSE
    IF s.lungime<>0 THEN
      BEGIN
        CreazăȘirVid(temp);
        Copiază(temp, s);
        {adaugă temp la sfârșitul lui u}
        u.coada^.urm:=temp.cap;
        u.coada:=temp.coada;
        u.lungime:=u.lungime+s.lungime
      END
    END; {ConcatSir}
-----

```

- Acest mod de reprezentare a șirurilor are unele **avantaje**.
 - (1) Permite o mare **flexibilitate** în gestionarea șirurilor.
 - (2) Înlătură **dezavantajul** lungimii fixe.

- Principalele **dezavantaje** rezidă în:
 - (1) **Parcurgerea secvențială** a șirului în vederea realizării accesului la o anumită poziție.
 - (2) **Risipa** de spațiu de memorie datorată prezenței pointerului asociat fiecărui caracter al șirului.

4.2.4. Comparație între metodele de implementare a șirurilor

- După cum s-a văzut, se utilizează trei maniere de implementare a șirurilor:
 - Implementarea bazată pe **tablouri**.
 - Implementarea bazată pe **pointeri**.
 - Implementarea bazată pe **tabele** care reunește parțial caracteristicile structurilor anterioare.
- Din punctul de vedere al paragrafului de față prezintă interes **primele două metode**, caracteristicile celei de-a treia putând fi ușor deduse din acestea.
- **Diferențele** marcante dintre cele două moduri de implementare a șirurilor și anume prin **tablouri** și **pointeri** își au de fapt originea în **caracterul fundamental** al **structurilor de date suport**. Astfel:
 - Șirul implementat ca **tablou** este o **structură de date statică** cu **acces direct**.
 - Șirul implementat prin **pointeri** este o **structură de date dinamică** cu **acces secvențial**, de fapt o **listă înlănțuită**.
- Aceste caracteristici fundamentale diferite influențează într-o manieră **decisivă** atât **proprietățile șirurilor** în fiecare din cele două moduri de reprezentare, cât și **algoritmii** care implementează **operatorii specifici**.
- **Implementarea șirurilor cu ajutorul tablourilor** oferă următoarele avantaje:
 - Posibilitatea realizării **directe** a accesului la caractere.
 - Accesul la o poziție precizată într-un șir este **imediat**, element care avantajează operațiile **CopiazăSubșir**, **ȘtergeȘir**, **InsereazăȘir**, **FurnizeazăCarȘir**.
 - Algoritmii specifici pentru **FurnizeazăCarȘir**, **AdaugăCarȘir** și **LungȘir** sunt rapizi ($O(1)$).
 - Restul operatorilor în marea lor majoritate au performanța $O(n)$.
 - Operatorul **PozițieSubșir** în implementarea bazată pe **căutare directă** are performanța $O(n * m)$ unde n este lungimea șirului iar m lungimea subșirului, (există pentru această reprezentare și metode mai performante).
 - **Dezavantajele** acestui mod de implementare sunt:
 - **Dimensiunea fixă** (limitată) a șirurilor.

- **Trunchierea** cauzată de inserția într-un șir complet.
- **Inserția** și **suprimarea** într-un șir presupun mișcări de elemente și un consum de timp proporțional cu lungimea șirului.
- **Implementarea șirurilor cu ajutorul pointerilor:**
 - Se bucură de **flexibilitatea** specifică acestei implementări.
 - Este grevată de proprietățile rezultate din **caracterul secvențial** al reprezentării.
 - Faza de inserție propriu-zisă și suprimare din cadrul operatorilor **InserazăȘir** și **ȘtergeȘir** se realizează în $O(I)$ unități de timp.
 - **Căutarea** poziției însă necesită practic $O(n)$ unități de timp.
 - De fapt **căutarea secvențială** a poziției este marele **dezavantaj** al acestei reprezentări care diminuează performanța mării majorități a operatorilor.
 - Cu toate că în fiecare moment se consumă spațiu de memorie numai pentru caracterele existente, totuși **risipa de memorie** datorată pointerilor este mare.
 - Performanța operatorilor **ConcatȘir** și **AdaugăCarȘir** crește în mod semnificativ dacă se utilizează pointerul "coadă".
 - Această implementare este indicat a fi utilizată atunci când:
 - Este foarte important ca șirurile să **nu** fie **trunchiate**.
 - Când se cunoaște cu probabilitate redusă lungimea maximă a șirurilor.
- Admițând trunchierea ca și un rău necesar, marea majoritate a limbajelor de programare care definesc și implementează în cadrul limbajului tipul `string` și operatorii aferenți utilizează implementarea șirurilor bazată pe **tablouri**.

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 **tehnicilor 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 tabelară (Table search)

- O căutare într-un tablou se numește și **căutare tabelară** ("**table search**"), cu deosebire în situația în care cheile sunt la rândul lor **obiecte structurate** spre exemplu **tablouri de numere** sau de **caractere**.

- Cazul mai des întâlnit este acela în care cheile tablourilor sunt **șiruri de caractere** sau **cuvinte**.
- În continuare, pentru obiectivele acestui paragraf, se definesc pentru o **structură șir**, în termeni de logica predicatelor, **relația de coincidență** respectiv **relația de ordonare** (lexicografică) a două șiruri x și y după cum urmează [4.3.1.a,b,c]:

/*Definirea tipului șir*/

typedef char tip_sir[m]; **/*[4.3.1.a]*/**

{Definirea relației de egalitate a două șiruri (coincidența)}

$(x=y) \leq (\mathbf{A_j : 0 \leq j < m : x_j = y_j})$ **{ [4.3.1.b] }**

{Definirea relației de ordonare a două șiruri}

$(x < y) \leq \mathbf{E_i : 0 \leq i < m : ((A_j : 0 \leq j < i : x_j = y_j) \ \& \ (x_i < y_i))}$ **{ [4.3.1.c] }**

$(x > y) \leq \mathbf{E_i : 0 \leq i < m : ((A_j : 0 \leq j < i : x_j = y_j) \ \& \ (x_i > y_i))}$

- Pentru a stabili **coincidența**, este necesară stabilirea **egalității tuturor caracterelor** corespunzătoare din șirurile comparate.
 - Presupunând că lungimea șirurilor este mică, (spre exemplu mai mică decât 30), în acest scop se poate utiliza **căutarea liniară**.
- În cele mai multe aplicații se lucrează cu șiruri de **lungime variabilă**, asociindu-se fiecărui șir o **informație suplimentară** referitoare la **lungimea** sa.
 - Utilizând tipul șir anterior definit, lungimea **nu** poate depăși valoarea maximă m .
 - Deși este limitativă, această schemă permite suficientă **flexibilitate** evitând alocarea dinamică a memoriei.
- Sunt utilizate în mod frecvent **două moduri de reprezentare a lungimii șirurilor**:
 - (1) Lungimea este specificată **implicit**, plasând pe ultima poziție a șirului (după ultimul caracter) un caracter prestabilit (**marker**) - de exemplu caracterul având codul 0 ($\text{CHR}(0)$).
 - Pentru aplicațiile ce urmează este important ca marker-ul să fie cel mai **mic** caracter al setului de caractere.
 - În limbajul C este chiar codul cu valoarea zero.
 - (2) Lungimea șirului este memorată în mod **explicit** ca și **prim element** al tabloului.
 - Astfel un șir are forma $s = s_0, s_1, s_2, \dots, s_{n-1}$ unde s_1, \dots, s_{n-1} sunt caracterele șirului iar s_0 memorează **lungimea șirului de caractere**.

- **Avantajul** acestei soluții: lungimea șirului este direct disponibilă.
- **Dezavantajul**: lungimea este limitată la valoarea maximă reprezentabilă pe unitatea de informație alocată unui caracter, de regulă un octet (255).
- În continuare, pentru precizarea **lungimii unui șir** se va utiliza modalitatea (1), respectiv cea bazată pe caracterul **marker**, definit drept caracterul al cărui cod are valoarea 0.
- În aceste condiții, **compararea** a două șiruri va lua forma [4.3.1.d]:

{Compararea liniară a două șiruri}

```
tip_sir x,y;                                     { [4.3.1.d] }

i=0;                                             {performanța O(n) }
while ((x[i]=y[i]) AND (x[i]<>chr(0))) i=i+1;
```

- Caracterul corespunzător codului 0 acționează ca și **fanion**.
- **Invariantul** buclei, adică condiția a cărei îndeplinire determină reluarea buclei (ciclului), este [4.3.1.e], iar **condiția de terminare** [4.3.1.f]:

{Invariantul buclei de comparare}

$A_j: 0 \leq j \leq i: x_j = y_j \neq \text{chr}(0)$ { [4.3.1.e] }

{Condiția de terminare a buclei de comparare}

$((x_i \neq y_i) \text{ OR } (x_i = y_i = \text{chr}(0))) \text{ \& } (A_j: 0 \leq j < i: x_j = y_j \neq \text{chr}(0))$ { [4.3.1.f] }

- **Condiția de terminare a buclei de comparare** stabilește:
 - (1) **Coincidența** celor două șiruri $x=y$ **dacă** la terminarea buclei $x_i=y_i=\text{chr}(0)$.
 - (2) **Inegalitatea** celor două șiruri $x<y$ ($x>y$) **dacă** la terminarea buclei $x_i<y_i$ ($x_i>y_i$).
- Se face precizarea că inegalitatea $x_i<y_i$ poate apare și în situația în care x **coincide** cu y dar x e mai **scurt** decât y .
 - În acest caz se obține $x_i \neq y_i$ dar $x_i = \text{chr}(0)$ (s-a ajuns la sfârșitul șirului x).
 - Întrucât în această situație x_i este cel mai mic caracter, aceasta echivalează cu $x_i < y_i$, deci $x < y$, ceea ce este corect, inclusiv din punctul de vedere al **ordonării lexicografice**.
- În acest context, **căutarea tabelară** este de fapt o **căutare încuibată** care constă:
 - (1) Dintr-o parcurgere a intrărilor unei table de șiruri.

- (2) Dintr-o secvență de comparații între componentele individuale ale șirului curent și cele ale șirului căutat, realizată în cadrul fiecărei intrări.
- Se consideră următoarele structuri de date [4.3.1.g], unde T este tabela de șiruri iar x : TipSir, argumentul căutat.

```

/*Căutarea tabelară - structuri de date*/
/*[4.3.1.g]*/
int n = /*numărul de șiruri din tabela de șiruri*/;
int m = /*numărul maxim de caractere dintr-un șir*/;
/*sfârșitul unui șir se marchează cu chr(0)*/

/*definire tip șir*/
typedef char Tip_Sir[m];
/*definire tabela de șiruri*/
typedef Tip_Sir Tip_Tabela[n];

Tip_Tabela T; /*instantiere tabelă de șiruri*/
Tip_Sir x; /*instantiere șir de căutat*/

```

- Presupunând că n este suficient de mare și că **tabela** este **ordonată alfabetic**, se poate utiliza **căutarea binară**.
- Utilizând algoritmul **căutării binare** și **compararea a două șiruri**, se poate redacta următoarea variantă de **căutare tabelară** [4.3.1.h].

```

/*Căutarea tabelară - varianta pseudocod*/

boolean CautareTabelara(Tip_Tabela T, Tip_Sir x, int * poz)

    int i,S,D,mijloc; /*[4.3.1.h]*/
    S=0; D=n;
    cat timp (S<D) /*căutare binară*/
    |   mijloc= (S+D)div2; i=0;
    |   cat timp ((T[mijloc,i]=x[i]) AND (x[i]<>chr(0))) i=i+1;
    |   daca (T[mijloc,i]<x[i]) S=mijloc+1;
    |   else D=mijloc;
    |   □ /*cât timp*/
    if (D<n) /*comparația finală de șiruri*/
    |   i= 0;
    |   cat timp ((T[D,i]=x[i]) AND (x[i]<>chr(0)))
    |   |   i=i+1;
    |   □ /*dacă*/
    *poz= D;
    return ((D<n) AND (T[D,i]=x[i]))
/*CautareTabelara*/

```

- Funcția **CautareTabelară** primește ca parametri de intrare tabela T și șirul de cautat x și returnează **true** respectiv **false** după cum căutarea este **reușită** sau **nu**.

- În cazul unei **căutări reușite**, în variabila de ieșire `poz` se returnează intrarea corespunzătoare din tabelă.
- În cazul unei **căutări nereușite** variabila de ieșire `poz` returnează o valoare nedeterminată.

4.3.2. Căutarea de șiruri directă

- O manieră frecvent întâlnită de căutare este așa numita **căutare de șiruri directă** ("string search").
- **Specificarea problemei:**
 - Se dă un tablou `s` de `n` caractere și un tablou `p` de `m` caractere, unde $0 < m < n$ declarate ca în secvența [4.3.2.a].
 - **Căutarea de șiruri directă** are drept scop stabilirea **primei apariții** a lui `p` în `s`.

*/*Căutarea de șiruri directă - structuri de date*/*

```
typedef unsigned int char;           /*[4.3.2.a]*/
char s[n];      /*șir original (text)*/
char p[m];      /*șir model de căutat (pattern)*/
```

- 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 **simplică** metodă de căutare este așa numita **căutare de șiruri directă** ("string search").
- **Rezultatul** unei astfel de căutări este un indice `i` care precizează apariția unei **coincidențe** de lungime `j` caractere între model și șir.
 - Acest lucru este formalizat de **predicatul** $P(i, j)$ [4.3.2.b].

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

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

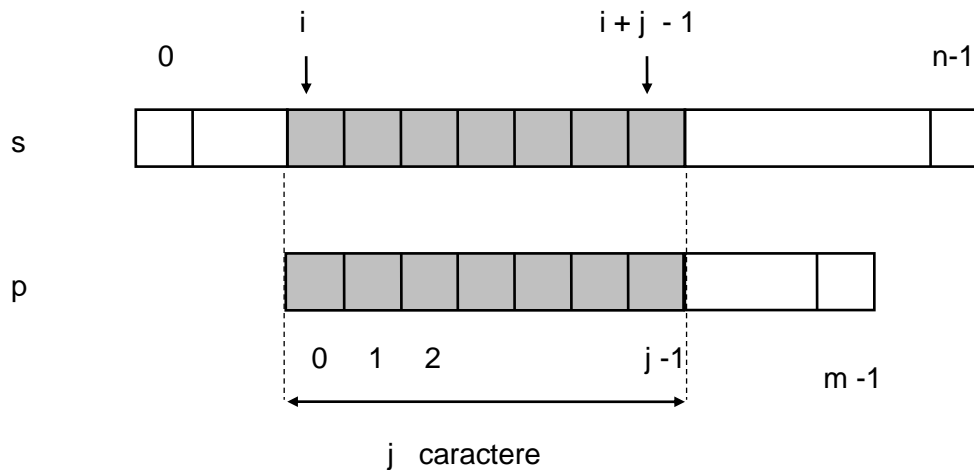


Fig.4.3.2.a. Reprezentarea grafică a predicatului $P(i, j)$

- 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].

$$Q(i) = \bigwedge_{k: 0 \leq k < i} \sim P(k, m)$$

[4.3.2.c]

- Specificarea problemei sugerează implementarea **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]:

{Căutarea de șiruri directă - Implementare a predicatului $Q(i)$ }

```

i=-1;                                     /*[4.3.2.d]*/
repetă
  i=i+1;
  gasit=P(i,m)
pană când (gasit OR (i=n-m))
□ /*repetă*/
  
```

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

/*Implementarea căutării de șiruri directe - varianta C*/

```
boolean cautare_directa(int* poz)           /*[4.3.2.e]*/
{
    int i;    /*i parcurge caracterele din șir*/
    int j;    /*j parcurge caracterele din model*/

    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;
    return (j==m);
} /*cautare_directa*/
```

- 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.

4.3.2.1 Analiza căutării de șiruri directe

- Algoritmul lucrează destul de **eficient** dacă se presupune că **nepotrivirea** în procesul de căutare apare după cel mult 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'.
 - 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.

4.3.3. Căutarea de șiruri Knuth-Morris-Pratt

- În anul 1970 Knuth, Morris și Pratt au inventat un algoritm de căutare în șiruri de caractere care necesită un **număr de comparații** de ordinul n chiar în cel mai **defavorabil** caz.
- Noul algoritm se bazează pe **observația** că avansul modelului în cadrul șirului cu o **singură** poziție la întâlnirea unei nepotriviri, așa cum se realizează în cazul **căutării directe**, pe lângă o eficiență scăzută, conduce la **pierderea** unor informații utile.
- Astfel după o **potrivire parțială** a modelului cu șirul:
 - Întrucât se cunoaște **parțial** șirul (până în punctul baleat);
 - Dacă avem cunoștințe **apriorice** asupra modelului obținute prin **precompilare**;
 - Le putem folosi pentru a **avansa mai rapid** în șir în procesul de căutare.
- Acest lucru este pus în evidență de **exemplul** următor în care:
 - Se caută modelul MARGINE în textul sursă dat.
 - Caracterele care se compară sunt cele subliniate.
 - După fiecare nepotrivire, modelul este deplasat de-a lungul **întregului** drum parcurs întrucât deplasări mai scurte **nu** ar putea conduce la o potrivire totală [4.3.3.a].

MAREA MARMARA SE MARGINESTE...

MARGINE

[4.3.3.a]

MARGINE

MARGINE

MARGINE

MARGINE

MARGINE

. . .

MARGINE

-
- Utilizând predicatele **P** și **Q**, **algoritmul KMP** poate fi formulat astfel [4.3.2.b]:

/*Căutarea de șiruri Knuth-Morris-Pratt*/

```
i=0; j=0;                                     /*[4.3.3.b]*/
cat timp ((j<m) AND (i<n))
| /*Q(i-j) && P(i-j,j)*/
| cat timp ((j>=0) AND (s[i]<>p[j])) j=d; /*deplasare*/
| i=i+1; j=j+1
| □ /*cat timp*/
```

- Formularea algoritmului este **aproape** completă, cu excepția factorului d care precizează **valoarea deplasării**.
- Se subliniază faptul că în continuare $Q(i-j)$ și $P(i-j, j)$ sunt **invarianții globali** ai procesului de căutare, la care se adaugă relațiile $0 < i < n$ și $0 < j < m$.
- Este important de subliniat faptul că, din rațiuni de claritate, predicatul P va fi ușor modificat: de aici înainte **nu** i va fi indicele care precizează **poziția primului caracter** al modelului în șir **ci** valoarea $i-j$.
 - Indicele i precizează **poziția la care a ajuns** procesul de căutare în șirul s . (fig.4.3.3.a).

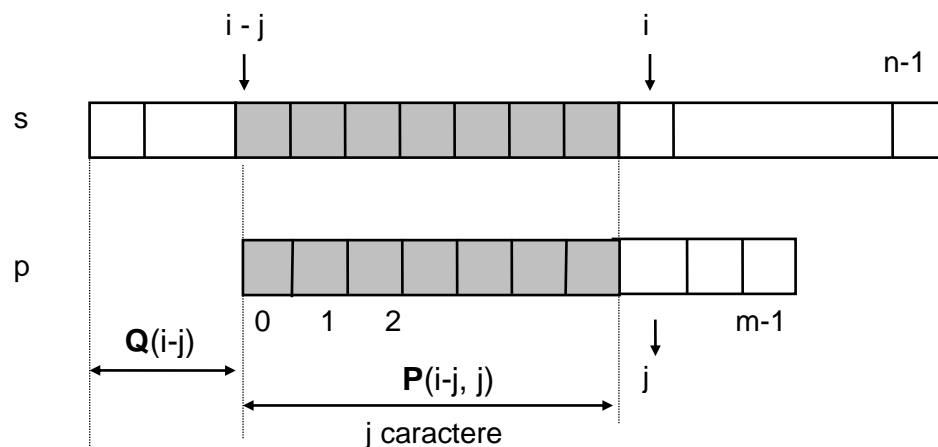


Fig. 4.3.3.a. Reprezentarea predicatelor P și Q în contextul căutării KMP

- Dacă algoritmul se termină deoarece $j=m$, termenul $P(i-j, j)$ al invariantului devine $P(i-m, m)$, ceea ce conform relației care-l definește pe P indică o **potrivire** între șir și model, începând de la poziția $i-m$ a șirului pe întreaga lungime m a modelului.
- Dacă la terminare $i=n$ și $j < m$, invariantul $Q(i)$ implică **absența** vreunei potriviri.
- În vederea determinării valorii lui d trebuie lămurit rolul atribuirii $j = d$.
- Considerând prin **definiție** că $d < j$ atunci atribuirea $j=d$ reprezintă o **deplasare a modelului** în cadrul șirului, spre **dreapta**, cu $j-d$ **poziții**.
 - Pentru a înțelege acest lucru trebuie remarcat faptul că **indicii** i și j sunt **întotdeauna aliniați** în procesul de căutare.
 - Este evident că este de dorit ca deplasarea să fie cât mai **lungă** și în consecință d cât mai **mic** posibil.
- Pentru a determina valoarea lui d se prezintă două exemple.

- **Exemplul 4.3.3.a.** Se consideră situația din figura 4.3.3.b.

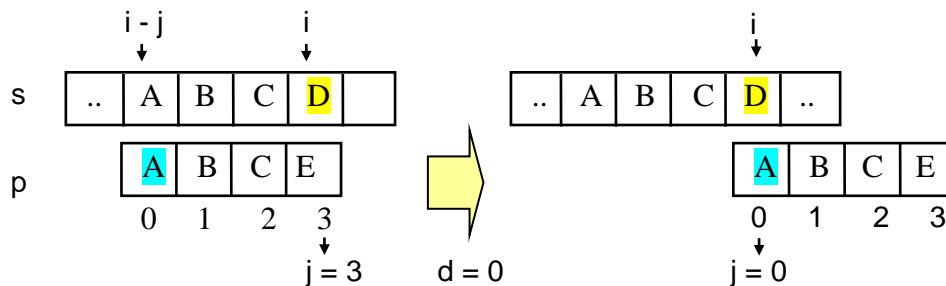


Fig. 4.3.3.b. Deplasarea modelului în cadrul șirului s . Cazul 1.

- Se observă că deplasarea se poate face peste 3 poziții întrucât nu mai pot apare alte potriviri. În aceste condiții $d=0$, iar atribuirea $j=d$ produce deplasarea cu $j-d=3$ poziții.

- **Exemplul 4.3.3.b.** Se consideră situația din figura 4.3.3.c.

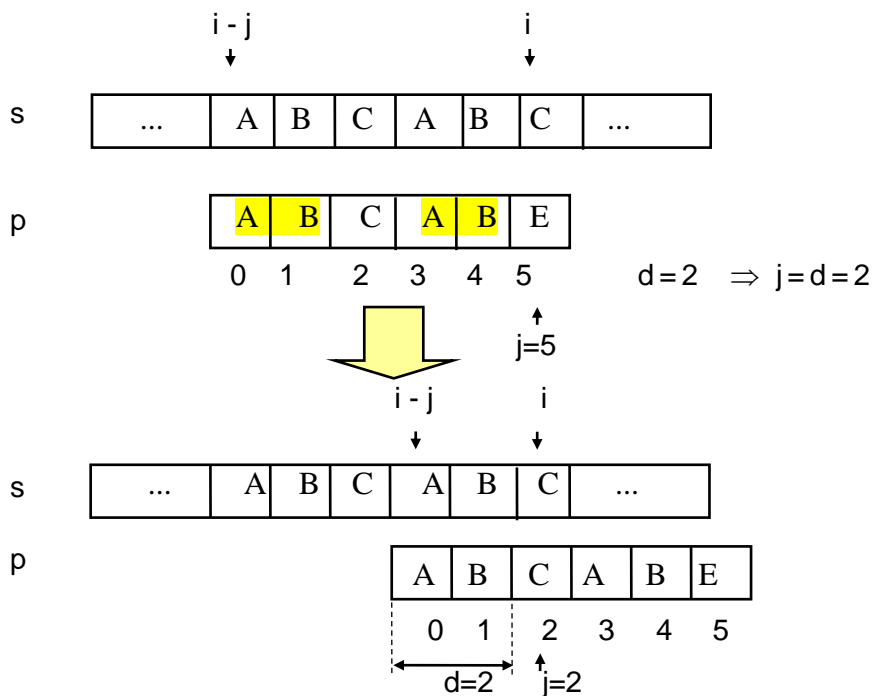


Fig. 4.3.3.c. Deplasarea modelului în cadrul șirului s . Cazul 2.

- În această situație deplasarea **nu** se poate face peste toată lungimea modelului întrucât mai există o posibilitate de potrivire începând de la $i-2$.

- Acest lucru se întâmplă întrucât în cadrul modelului există 2 **subsecvențe identice** 'AB' de lungime 2.
- Se notează cu d lungimea acestei subsecvențe.
- După cum se observă din figura 4.3.3.c, deplasarea modelului la dreapta se poate face **numai** peste $j-d$ ($5-2=3$) poziții deoarece **primele** d poziții ale modelului **coincid** cu d poziții ale modelului care preced nepotrivirea (indicele j) și implicit cu cele d caractere ale textului, care **preced** indicele i (care indică neconcordanța).

-
- **În concluzie:**
 - **Inițial** există o potrivire între șirul s și modelul p începând de la poziția $i-j$ pe lungime j adică este valabil invariantul $P(i-j, j) \ \& \ Q(i-j)$.
 - **După** efectuarea deplasării avem din nou o potrivire la $i-d$ de lungime d adică este valabil invariantul $P(i-d, d) \ \& \ Q(i-d)$.
 - Această situație rezultă din observația că **în model** există **două subsecvențe identice** de lungime d una la începutul modelului alta de la poziția $j-d$ la $j-1$.
 - **Formal** acest lucru este precizat de relația [4.3.3.c].

$$p_0 \dots p_{d-1} = p_{j-d} \dots p_{j-1} \qquad [4.3.3.c]$$

- Pentru ca lucrurile să se desfășoare **corect** este necesar ca lungimea lui d să fie **maximă**, adică să avem **cea mai lungă subsecvență** care îndeplinește condițiile precizate.
 - **Rezultatul esențial** este acela că **valoarea** lui d este determinată **exclusiv** de **structura modelului** și **nu** depinde de textul propriu-zis.
 - Din cele expuse rezultă că pentru a-l determina pe d se poate utiliza următorul **algoritm**:
 - Pentru fiecare poziție j din cadrul modelului.
 - Se caută în model cel mai mare d , adică cea mai lungă secvență de caractere a modelului care precede imediat poziția j și care se potrivește ca un număr egal de caractere de la începutul modelului.
 - Se notează valoarea d pentru un anumit j cu d_j .
 - Întrucât valorile d_j depind **numai și numai** de **model**, **înaintea** căutării propriu-zise poate fi construit un **tablou** d cu aceste valori, printr-un proces de **precompilare a modelului**.
 - Desigur acest efort merită să fie depus dacă $m \ll n$.
 - În continuare se prezintă trei **situații particulare** de determinare a valorii deplasamentului d .
-

- **Exemplul 4.3.3.c.** Se consideră situația din figura 4.3.3.d.

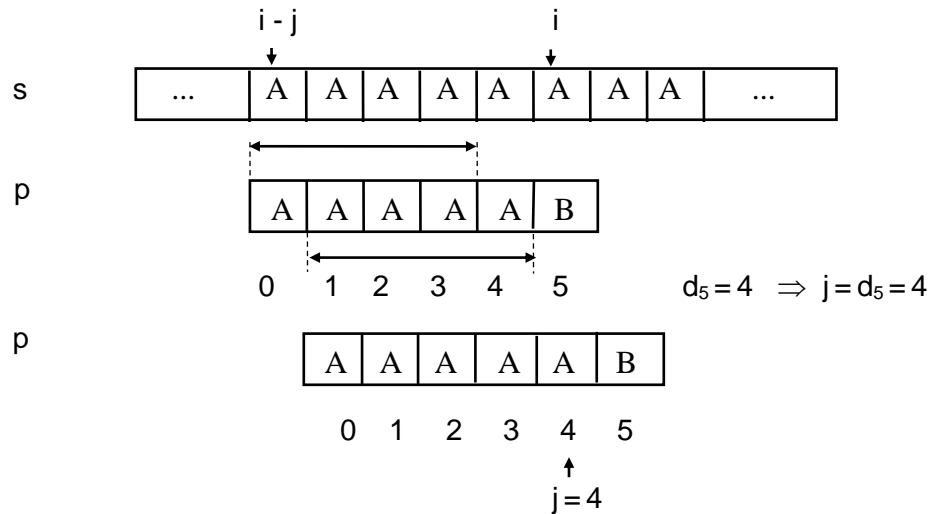


Fig. 4.3.3.d. Determinarea deplasamentului d . Cazul 1.

-
- **Exemplul 4.3.3.d.** Se consideră situația din figura 4.3.3.e.

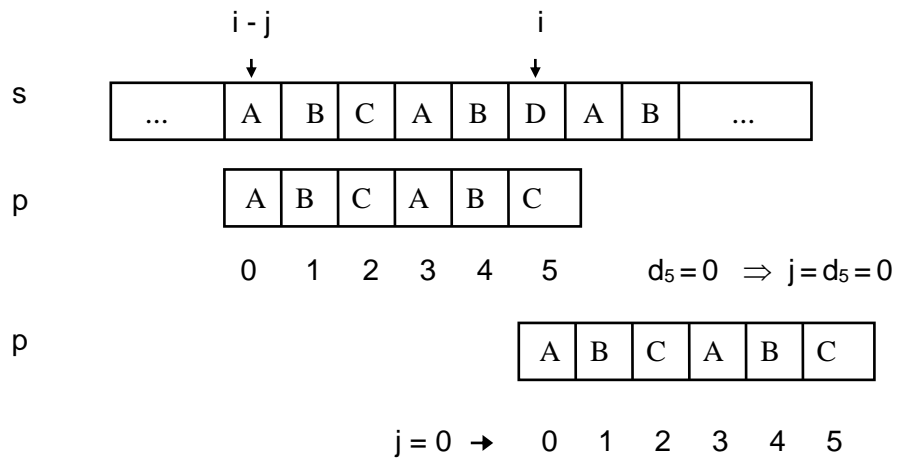
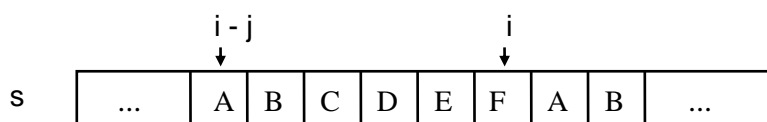


Fig. 4.3.3.e. Determinarea deplasamentului d . Cazul 2.

-
- **Exemplul 4.3.3.e.** Se consideră situația din figura 4.3.3.f.



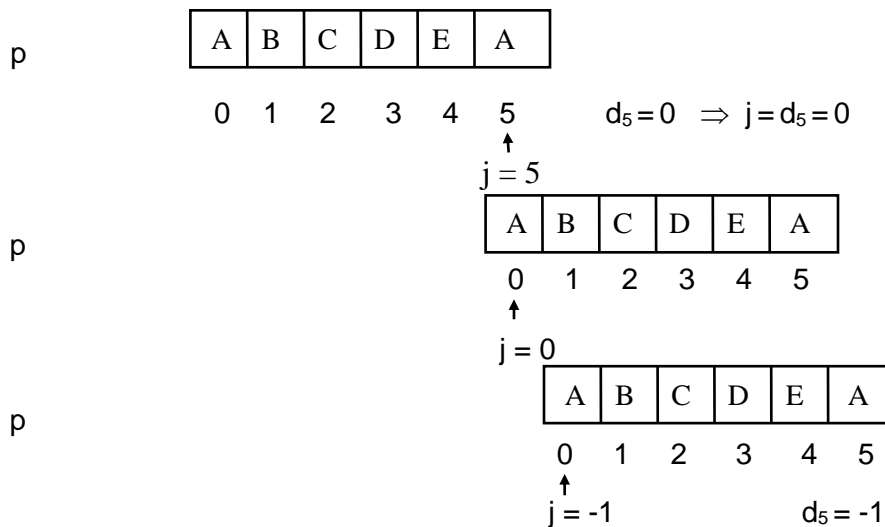


Fig. 4.3.3.f. Determinarea deplasamentului d . Cazul 3.

- Ultimul exemplu sugerează că se poate merge mai departe cu deplasările și că în loc să se realizeze deplasarea peste 5 poziții, în această situație se poate face peste întregul model, deci $d_5 = -1$ iar deplasarea se face peste $5 - (-1) = 6$ poziții.
 - Acest lucru este posibil deoarece **primul** caracter al modelului este **identic** cu **ultimul** său caracter și **diferit** de caracterul i al șirului.
 - Întrucât s-a constatat **necoincidența** pentru ultimul caracter al modelului, rezultă că **nu** poate exista o coincidență nici în cazul primului caracter al acestuia (identic cu ultimul), deci deplasarea se poate realiza inclusiv **peste** acesta.
- În aceste condiții, calculul lui d_j care presupune căutarea celor mai lungi secvențe care se potrivesc în baza relației [4.3.3.c], poate fi completat după cum urmează.
 - Dacă se constată că $d_j = 0$ și $p_0 = p_j$, atunci se poate face $d_j = -1$ indicând **deplasarea integrală** a modelului față de poziția sa curentă în cadrul șirului s .
- În figura 4.3.3.g. apare un **tablou demonstrativ** în care pentru mai multe șiruri model p se precizează structura tabloului d asociat adică valorile d_j corespunzătoare caracterelor modelului rezultate în urma precompilării.
- În stânga tabelului apare modelul p iar în dreapta tabelul d asociat.
- Programul care implementează **căutarea de șiruri Knuth-Morrison-Pratt** apare în secvența [4.3.3.d].

p							d						
0	1	2	3	4	5	6	0	1	2	3	4	5	6
A							-1						
A	A						-1	-1					

A	A	A	A	A	B		-1	-1	-1	-1	-1	4	
A	B	C	A	B	C		-1	0	0	-1	0	0	
A	B	C	A	B	C	D	-1	0	0	-1	0	0	3
A	B	C	A	B	D		-1	0	0	-1	0	2	
A	B	C	D	E	F		-1	0	0	0	0	0	
A	B	C	D	E	A		-1	0	0	0	0	-1	
M	A	R	G	I	N	E	-1	0	0	0	0	0	0

Fig. 4.3.3.g. Exemple de precompilare a unor modele

```

/*Căutarea de șiruri Knuth-Morrison-Pratt - pseudocod*/

int mmax=/*lungime maximă model*/;           /*[4.3.3.d]*/
int nmax=/*lungime maximă șir sursă*/;

int m          /*lungime model*/;
int n          /*lungime șir*/;
char  p[mmax]; /*model*/
char  s[nmax]; /*șir*/
int    d[mmax]; /*tabela de deplasări*/

boolean CautareKMP(int * poz)
    int i,j,k;

[1]  *citire șir în s          /*n = lungime curentă șir*/
[2]  *citire model în p       /*m = lungime curentă model*/
[3]  j=0; k=-1; d[0]=-1;      /*precompilare model*/
      cat timp(j<m-1)
      |  cat timp((k>=0)AND(p[j]<>p[k])) k=d[k];
      |  j=j+1; k=k+1;
      |  daca(p[j]=p[k]) d[j]=d[k];
      |  altfel d[j]=k;
      |  □ /*cât timp*/
[4]  i=0; j=0;                /*căutare model*/
      cat timp((j<m)AND(i<n))
      |  cat timp((j>=0)AND(s[i]<>p[j])) j=d[j];
      |  j=j+1;
      |  i=i+1;
      |  □ /*cât timp*/
      *poz=i-m;
      return (j==m);
/*CautareKMP*/

```

- Programul KMP constă din 4 părți:
 - [1] Prima parte realizează **citirea șirului s** în care se face căutarea.
 - [2] A doua parte realizează **citirea modelului p**.

- [3] A treia parte **precompilează modelul** și calculează valorile d_j .
- [4] Cea de-a patra parte **implementează căutarea** propriu-zisă.
- Se precizează că partea a treia a algoritmului (precompilarea modelului) fiind practic o căutare de șiruri, se implementează tot ca și o **căutare KMP** utilizând porțiunea deja determinată a tabloului precompilat d .

4.3.3.1 Analiza căutării de șiruri Knuth-Morrison-Pratt

- Analiza exactă a performanței căutării de șiruri Knuth-Morrison-Pratt este asemenea algoritmului, foarte sofisticată.
- Inventatorii demonstrează că **numărul de comparații** de caractere este de ordinul $n+m$.
 - Aceasta reprezintă o îmbunătățire substanțială față de $m*n$.
- De asemenea se subliniază faptul că pointerul i care baleează șirul **nu** merge niciodată înapoi.
 - Acest lucru diferențiază căutarea KMP de căutarea directă, unde după o potrivire parțială se reia căutarea cu modelul deplasat cu o poziție, chiar dacă o parte din caracterele șirului au fost deja parcurse.
 - Acest avantaj permite aplicarea acestei metode și în cazul unor **prelucrări secvențiale**.

4.3.4. Căutarea de șiruri Boyer-Moore

- Metoda ingenioasă de căutare KMP conduce la beneficii **numai** dacă **nepotrivirea** dintre șir și model a fost **precedată** de o **potrivire parțială** de o anumită lungime.
 - **Numai** în acest caz deplasarea modelului se realizează peste mai **mult** de o poziție.
- Din păcate această situație în realitate este mai degrabă **excepția** decât **regula**; potrivirile apar mult mai rar ca și nepotrivirile.
 - În consecință, **beneficiile** acestei metode **sunt reduse** în marea majoritate a căutărilor normale de texte.
- Metoda de căutare inventată în 1975 de **R.S. Boyer** și **J.S. Moore** îmbunătățește performanța atât pentru situația cea mai defavorabilă cât și în general.
- Ideea de căutării este clasică: modelul avansează de-a lungul șirului și se compară identitatea caracterelor corespunzătoare dintre model și șir. În caz de nereușită, modelul avansează de-a lungul șirului.

- **Căutarea BM**, după cum mai este numită, este bazată însă pe ideea **neobișnuită** de a începe compararea caracterelor între șir și model de la **sfârșitul modelului** și **nu** de la **începutul** acestuia.
- Ca și în cazul metodei KMP, modelul este **precompilat** anterior într-un tablou d .
- **Precompilarea** presupune următorii pași:
 - (1) Pentru **fiecare caracter x care apare în model**, 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.).
 - (2) Valoarea d_x se trece în tabloul d în poziția corespunzătoare caracterului x .
 - (3) Pentru toate celelalte caractere ale setului de caractere, **care nu apar în model** d_x se face egal cu **lungimea totală** a modelului.
 - (4) Pentru **ultimul caracter al modelului**, d_x se face de asemenea egal cu **lungimea totală** a modelului.

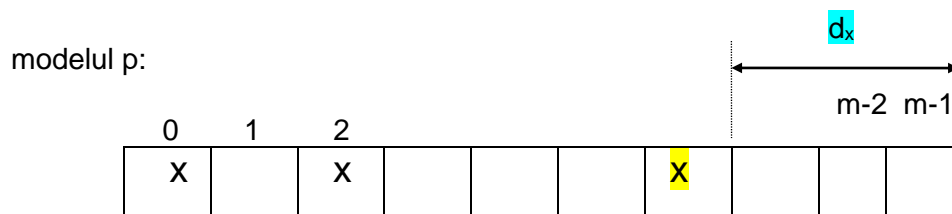


Fig. 4.3.4.a. Determinarea valorii d_x corespunzătoare caracterului x al modelului

- În continuare, se presupune că în procesul de **comparare de la dreapta la stânga** 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.).

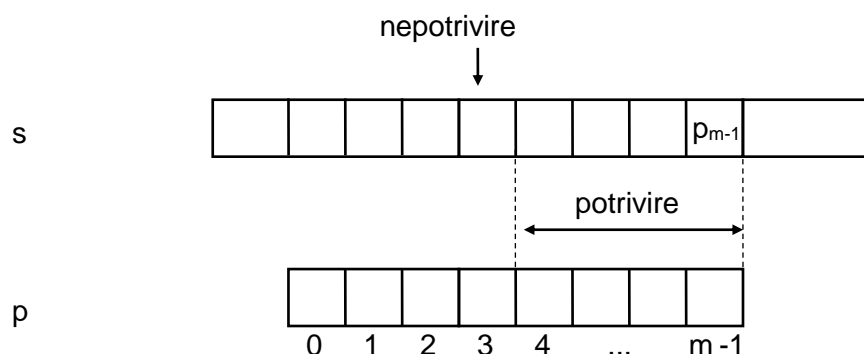
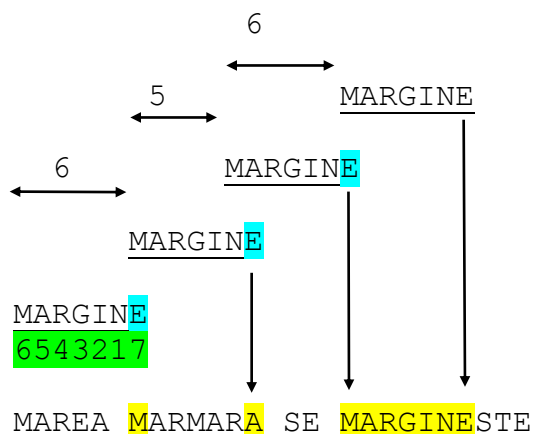


Fig. 4.3.4.b. Comparare șir-model pentru determinarea lui p_{m-1}

- 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.



[4.3.4.a]

- Deoarece compararea individuală a caracterelor se realizează de la dreapta spre stânga, este mai **convenabilă** următoarea reformulare a predicatelor P și Q [4.3.4.b].

$$P(i, j) = A_k : j \leq k < m : s_{i-j+k} = p_k$$

$$Q(i) = A_k : 0 \leq k < i : \sim P(k, 0)$$

[4.3.4.b]

- Interpretarea grafică a a noii formulări a predicatului $P(i, j)$ în cauză apare în figura 4.3.4.c.

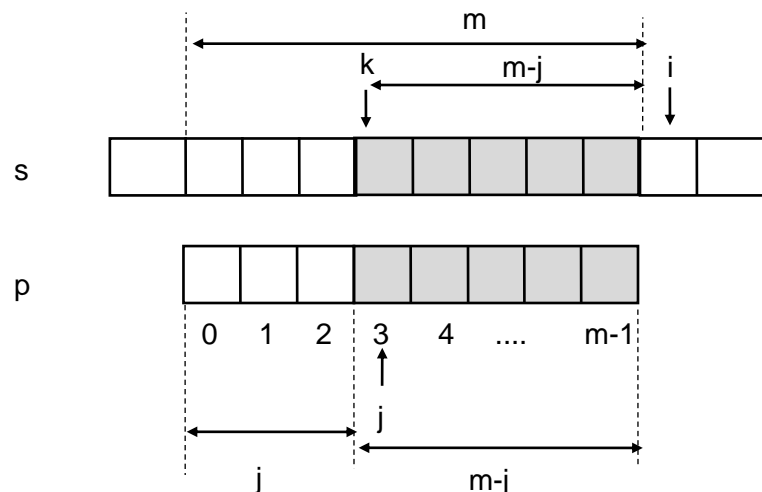


Fig. 4.3.4.c. Interpretarea predicatelor P și Q în căutarea BM

- În această interpretare indicele i indică poziția următoare celei de la care se demarează căutarea curentă (spre stânga).
 - După cum se observă, există o **potrivire** începând de la i spre **stânga** cu ultimele $m-j$ caractere ale modelului, adică începând de la sfârșitul modelului, până la poziția j (spre stânga).
 - Când $j=0$, avem o potrivire de la $i-m$ cu $m-0$ caractere, deci cu modelul integral (de la poziția 0 a modelului).
- Aceste predicate sunt utilizate în formularea următoare a algoritmului BM [4.3.4.c]:

```
/*căutare Boyer-Moore - varianta C*/
```

```

i=m; j=m;                                     /*[4.3.4.c]*/
while ( (j>0) && (i<n) )
{
    /*Q(i-m)*/
    j=m; k=i;
    /*comparare șir-model*/
    while ( (j>0) && (s[k-1]==p[j-1]) )
    {
        /*P(k-j, j) && (k-j=i-m)*/
        k=k-1;
        j=j-1;
    } /*while*/
    if (j>0) i=i+d[s[i-1]]; /*deplasare model în șir*/
} /*while*/
/*cautare Boyer-Moore*/

```

- Indicii implicați satisfac următoarele relații: $0 < j < m$, $0 < i$ și $k < n$.
 - **Terminarea** algoritmului cu $j=0$ presupune că $P(k-j, j)$ devine $P(k, 0)$ ceea ce indică o **potrivire** începând de la poziția k din șir, spre dreapta, de m caractere, unde $k=i-m$.
 - **Terminarea** algoritmului cu $j>0$ implică $i=n$.
 - În acest caz $Q(i-m)$ devine $Q(n-m)$ indicând **absența potrivirii**.
- Programul următor [4.3.4.d] implementează **strategia Boyer-Moore** într-un context similar căutării KMP.

```
/* Căutarea Boyer-Moore -varianta C */
```

```

typedef unsigned char boolean;                /*[4.3.4.c]*/
#define true (1)
#define false (0)

enum { mmax = 100 /*lungime maximă model*/,
      nmax = 200 } /*lungime maximă șir sursă*/;
int m; /*lungime model*/

```



```

int n;                                /*lungime şir*/;
char p[mmax];                         /*model*/
char s[nmax];                         /*şir*/
int d[256];                           /*tabela de deplasări*/

boolean cautare_bm(int* poz)
{
    int i,j,k;

    *citire sir;    /*n este lungimea curentă a şirului*/
    *citire model;  /*m este lungimea curentă a modelului*/
    /*inițializare tabelă de deplasări*/
    for(i=0;i<=255;i++) d[i]=m;
    /*precompilare model, construcție tabelă de deplasări*/
    for(j=0;j<=m-2;j++) d[p[j]]=m-j-1;
    /*căutare model*/
    i=m; j=m;
    while((j>0) && (i<=n))
    {
        j=m;
        k=i;
        while((j>0) && (s[k-1]==p[j-1]))
        {
            k=k-1;
            j=j-1;
        } /*while*/
        if(j>0) i=i+d[s[i-1]];
    } /*while*/
    *poz=i-m; /*poz=k*/
    return j==0;
} /*cautare_bm*/
/*-----*/

```

4.3.4.1 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 şir peste un caracter diferit de cele ale modelului, **numărul de comparații** este n/m .
- Autorii indică anumite direcții de **îmbunătățire a performanțelor** algoritmului.
 - Una din ele este aceea de a **combina** strategia **BM** care realizează o deplasare substanțială în prezența unei nepotriviri cu strategia **KMP** care permite o deplasare mai substanțială după detecția unei potriviri (parțiale).
 - Această metodă necesită **două** tablouri precalculate:
 - d_1 - tabloul specific căutării **BM**.
 - d_2 - tabloul corespunzător căutării **KMP**.

- Toate acestea conduc la **complicarea** algoritmului și la **creșterea regiei** sale prin precompilarea tablourilor.
- De fapt în multe cazuri trebuie analizată oportunitatea implementării unor extensii sofisticate care deși **rezolvă anumite situații punctuale**, pot avea drept consecință **deteriorarea performanțelor de ansamblu** prin **creșterea excesivă a regiei** algoritmului implicat.