

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 **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, 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.

Notății: *s, sub, u* - siruri;
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**;

CopiazaSubsir(*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*; eventual
 trunchiere la depășirea lungimii maxime;

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 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 InseereazăȘir utilizând operatorii CreazăȘirVid și AdaugăCar}

```
PROCEDURE InseereazaSir(subsir: TipSir; VAR sursa: TipSir;
                        p: TipIndice);
{inserează "subșir" în "sursă" între pozițiile p și p+1}
VAR rezultat: TipSir; i: TipIndice;
BEGIN
  IF (p<1) OR (p>LungSir(sursa)) THEN
    *eroare(pozitie ilegală in inserție)
  ELSE
    BEGIN
      CreazaSirVid(rezultat);
      FOR i:=1 TO p-1 DO
        AdaugaCarSir(rezultat, FurnizeazaCarSir(sursa,i));
      FOR i:=1 TO LungSir(subsir) DO
```

[4.1.b]

```

        AdaugaCarSir(rezultat, FurnizeazaCarSir(subsir,i));
FOR i:=p TO LungSir(sursa) DO
    AdaugaCarSir(rezultat, FurnizeazaCarSir(sursa,i));
    CreazaSirVid(sursa);
FOR i:=1 TO LungSir(rezultat) DO
    AdaugaCarSir(sursa, FurnizeazaCarSir(rezultat,i));
END {ELSE}
END; {InsereazaSir}
-----
/*Implementarea operatorului Inserează_Şir utilizând
operatorii Crează_Sir_Vid, Adaugă_Car_Sir şi
Furnizeaza_Car_Sir */

void insereaza_sir(tip_sir subsir, tip_sir* sursa,
                  tip_indice p)
/*inserează subsirul în sursa între poziţiile p şi p+1*/
{
    tip_sir rezultat; tip_indice i;

    if ((p<1)|| (p>lung_sir(*sursa))) ;
        /*eroare (poziţie ilegală în inserţie)*/
    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 şir

- Sunt cunoscute **două tehnici majore** de implementare a şirurilor:
 - (1) Implementarea bazată pe **tablouri**
 - (2) 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.
- 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 secvența [4.2.1.a].

{Exemplu de implementare a TDA Sir utilizând tablouri -
structuri de date - varianta Pascal}

```
CONST LungimeMax = ...;
TYPE TipLungime = 0..LungimeMax;
    TipIndice = 1..LungimeMax;           [4.2.1.a]
    TipSir = RECORD
        lungime: TipLungime;
        sir: ARRAY[TipIndice] OF char
    END;
VAR s: TipSir;
```

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

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

typedef unsigned char tip_indice;           /*[4.2.1.a]*/

typedef struct tip_sir {
    tip_lungime lungime;
    char sir[lungime_max];
} tip_sir;
tip_sir 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 **santinelă de sfârșit (marker)**.
 - În această situație operatorul `LungimeSir` va trebui să **contorizeze** într-o manieră liniară caracterele până la detectarea markerului.

- Din acest motiv este preferabil ca **lungimea șirului** 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 `LungimeMaximă`, 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 `ȘirComplet`.
- În accepțiunea modelului anterior prezentat, în continuare se prezintă o implementare a operatorilor primitivi [4.2.1.b,c,d,e] variantele Pascal și C

{Exemplu de implementare a TDA Șir utilizând tablouri - varianta Pascal}

```
PROCEDURE CreazaSirVid(VAR s: TipSir);           {O(1)}
BEGIN
    s.lungime:= 0                                [4.2.1.b]
END; {CreazaSirVid}
```

```
FUNCTION LungSir(VAR s: TipSir): TipLungime;     {O(1)}
BEGIN
    LungSir:= s.lungime                          [4.2.1.c]
END; {LungSir}
```

```
FUNCTION FurnizeazaCar(s: TipSir, poz: TipIndice): char; {O(1)}
BEGIN
    IF (poz<1) OR (poz>s.lungime) THEN
        BEGIN                                     [4.2.1.d]
            *eroare(pozitie incorecta);
            FurnizeazaCar:= chr(0) {caracterul vid}
        END
    ELSE
        FurnizeazaCar:= s.sir[poz]
    END; {FurnizeazaCar}
```

```
PROCEDURE AdaugaCarSir(VAR s: TipSir; c: char); {O(1)}
BEGIN
    IF s.lungime=LungimeMax THEN
        *eroare(se depășește lungimea maximă a șirului)
    ELSE
        BEGIN
            s.lungime:= s.lungime+1;              [4.2.1.e]
            s.sir[s.lungime]:= c
        END
    END; {AdaugaCarSir}
```

/* Exemplu de implementare a TDA Șir utilizând tablouri - varianta C*/

```

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

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

-----
char furnizeaza_car(tip_sir s,tip_indice poz)                 /*O(1)*/
{
    char furnizeaza_car_result;
    if ((poz<1)|| (poz>s.lungime))
    {
        /*eroare(pozitie incorecta);*/
        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)                         /*O(1)*/
{
    if (s->lungime==lungime_max) ;
    /*eroare(se depaseste lungimea maxima a sirului)*/
    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 **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}

PROCEDURE CopiazaSubsir(VAR u,s: TipSir, poz, lung:

TipLungime); {O(n)}

{copiază un subșir din s în șirul u}

VAR indexSursa, indexCopiere: TipLungime;

BEGIN

IF (poz<1) OR (poz>s.lungime) THEN

u.lungime:= 0

ELSE

BEGIN

[4.2.1.f]

indexSursa:= poz-1;

indexCopiere:= 0;

WHILE (indexSursa<s.lungime) AND

(indexCopiere<u.lungime) DO BEGIN

indexSursa:= indexSursa+1;

indexCopiere:= indexCopiere+1;

u.sir[indexCopiere]:= s.sir[indexSursa]

END; {WHILE}

u.lungime:= indexCopiere

END {ELSE}

END; {CopiazăSubsir}

/*Implementarea operatorului CopiazăSubsir*/ /*O(n)*/

{copiază un subșir din s în șirul u}

void copiaza_subsir(tip_sir* u, tip_sir* s, tip_lungime poz,
tip_lungime lung)

{

tip_lungime index_sursa, index_copiere;

if ((poz<1)|| (poz>s->lungime))

u->lungime= 0;

else

{

/*[4.2.1.f]*/

index_sursa= poz-1;

index_copiere= 0;

while ((index_sursa<s->lungime) &&

(index_copiere<u->lungime)) {

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;

}

}

-
- Se observă faptul că în interiorul buclei **WHILE** există 3 atribuiri, care pentru o lungime lung a sub șirului determină $3 \cdot \text{lung} + 3$ atribuiri.
 - Se observă de asemenea că procedura CopiazăSubsir este liniară în raport cu lungimea lung a subșirului.

- Un alt exemplu îl reprezintă implementarea funcției PozițieSubșir [4.2.1.g].

{Implementarea operatorului PozițieSubșir}

```

FUNCTION PozitieSubsir(VAR sub,s: tipSir): TipLungime;

{determină poziția lui "sub" în cadrul șirului "s" (prima
aparitie)}

VAR mark, {index pentru punctul de start al unei comparații}
    indexSub,          {index subșir}
    indexSir: TipIndex; {index șir}
    poz: TipLungime;    {poziția lui sub în s}
    stop: boolean; {devine adevărat când elementele
                    corespunzătoare din s și sub nu sunt egale}

BEGIN
    indexSir:= 1;
    poz:= 0;
    REPEAT
        indexSub:= 1;
        mark:= indexSir;
        stop:= false;
        WHILE (NOT stop) AND (indexSir<=s.lungime) AND
            (indexSub<=sub.lungime) DO
            IF s.sir[indexSir]=sub.sir[indexSub] THEN
                BEGIN
                    indexSir:= indexSir+1;
                    indexSub:= indexSub+1
                                [4.2.1.g]
                END
            ELSE
                stop:= true; {WHILE}
            IF indexSub>sub.lungime THEN
                poz:= mark {potrivire}
            ELSE
                indexSir:= mark+1 {nepotrivire => avans model}
            UNTIL (poz>0) OR (indexSir+sub.lungime-1>s.lungime);
    PozitieSubsir:= poz
END; {PozitieSubsir}

```

/*Implementarea operatorului PozițieSubșir*/

```

tip_lungime pozitie_subsir(tip_sir* sub,tip_sir* s)
{
    tip_index mark, /*index pentru punctul de start al unei
                    comparații*/
    index_sub,      /*index subsir*/
    index_sir;      /*index sir*/
    tip_lungime poz; /*poziția lui sub în s*/
    boolean stop;   /*devine adevărat când elementele
                    corespunzătoare din s și sub nu sunt egale*/

    tip_lungime pozitie_subsir_result;
    index_sir= 1;
    poz= 0;

```

```

do {
    index_sub= 1;
    mark= index_sir;
    stop= false;
    while ((!stop)&&(index_sir<=s->lungime)&&
           (index_sub<=sub->lungime))
        if (s->sir[index_sir-1]==sub->sir[index_sub-1])
        {
            index_sir= index_sir+1;
            index_sub= index_sub+1;          /*[4.2.1.g]*/
        }
        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));
pozitie_subsir_result= poz;
return pozitie_subsir_result;
}

```

- Complexitatea funcției **PozițieSubșir**(sub,s) este $O(lungime * s.lungime)$ 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.
- 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 inefficient 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.

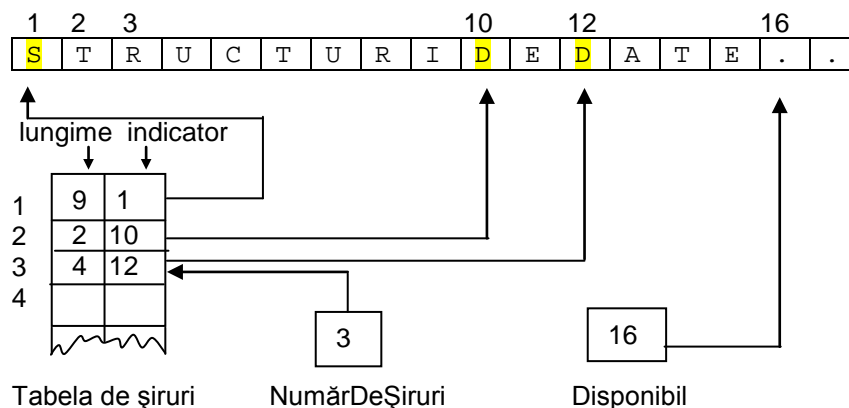


Fig.4.2.2.a. Modelul tabelii de șiruri

- Un exemplu de implementare în limbajul PASCAL 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 tabelii de șiruri.

{Implementarea șirurilor utilizând tabele - structuri de date}

```
CONST LungimeHeap = {număr maxim de caractere};
      LungimeTabela = {număr maxim de șiruri};
```

```
TYPE ElementTabela = RECORD
    lungime,
```

```

        indicator: 1..LungimeHeap                                [4.2.2.a]
    END;
    TipSir = 1..LungimeTabela;

VAR TabelaDeSiruri: ARRAY[1..LungimeTabela] OF
                                                ElementTabela;
    Heap: ARRAY[1..LungimeHeap] OF char;
    Disponibil: 0..LungimeHeap; {primul caracter liber}
    NumarDeSiruri: 0..LungimeTabela;

```

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

- Deoarece adăugarea unui caracter produce depăşirea şirului în cauză în dauna şirului următor din spaţiul de lucru, funcţia **AdaugăCarŞir** trebuie:
 - [1] Să creeze o nouă copie a şirului în zona disponibilă a heap-ului.
 - [2] Să recopieze şirul în noua locaţie.
 - [3] Să adauge noul caracter.
 - [4] Să reactualizeze tabela de şiruri şi variabila Disponibil [4.2.2.b].
-

{Implementarea operatorului **AdaugăCarŞir**}

```

PROCEDURE AdaugaCarSir(VAR s: TipSir; c: char); {O(n)}
    {adaugă caracterul c la sfârşitul sirului s}
    VAR i, lungimeVeche, indicatorVechi: integer;

    BEGIN
        IF (s<1)OR(s>NumarDeSiruri) THEN
            *eroare(referinţă invalidă de şir)
        ELSE
            BEGIN
                [1]   lungimeVeche:= TabelaDeSiruri[s].lungime;
                     indicatorVechi:= TabelaDeSiruri[s].indicator;
                [2]   FOR i:=0 TO lungimeVeche-1 DO
                     Heap[Disponibil+i]:=Heap[indicatorVechi+i];
                [3]   Heap[Disponibil+lungimeVeche]:= c; {se adauga c}
                [4]   TabelaDeSiruri[s].indicator:= Disponibil;
                     TabelaDeSiruri[s].lungime:= lungimeVeche+1;
                     Disponibil:= Disponibil+lungimeVeche+1
            END {ELSE}
        END; {AdaugaCarSir}

```

- Se observă că în această implementare **AdaugăCarŞir** necesită $O(n)$ unităţi de timp pentru un şir de n caractere faţă 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 înlanțuite**.
 - 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}

```
TYPE PointerCelula = ^Celula;  
    Celula = RECORD  
        ch: char;  
        urm: PointerCelula  
    END;  
    TipSir = RECORD  
        lungime: integer;  
        cap,coada: PointerCelula  
    END;
```

[4.2.3.a]

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

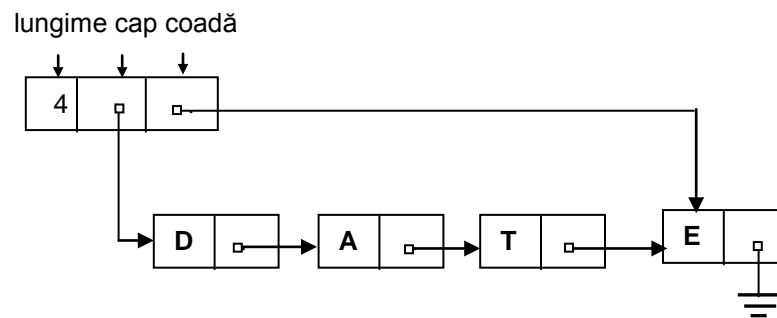


Fig.4.2.3.a. Implementarea șirurilor cu ajutorul pointerilor

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

{Implementarea șirurilor cu ajutorul pointerilor}

```
PROCEDURE CreazaSirVid(VAR s: TipSir);    {O(1)}  
BEGIN  
    s.lungime:= 0;  
    s.cap:= nil;  
    s:coada:= nil  
END; {CreazaSirVid}
```

[4.2.3.b]

```
FUNCTION LungSir(s: TipSir): integer;    {O(1)}  
BEGIN  
    LungSir:= s.lungime  
END; {LungSir}
```

[4.2.3.c]

```
FUNCTION FurnizeazaCarSir(s:TipSir; poz:integer): char;  
    VAR contor: integer;    {O(n)}  
        p: PointerCelula;  
BEGIN  
    IF (poz<1) AND (poz>s.lungime) THEN  
        BEGIN  
            *eroare(index eronat)
```

[4.2.3.d]

```

        FurnizeazaCarSir:= chr(0) {caracterul nul}
    END
ELSE
    BEGIN
        p:= s.cap;
        FOR contor:= 1 TO poz-1 DO p:= p^.urm;
        FurnizeazaCarSir:= p^.ch
    END
END; {FurnizeazaCarSir}
-----
PROCEDURE AdaugaCarSir(var s:TipSir; c:char); {O(1)}
BEGIN
    IF s.lungime=0 THEN {primul caracter al sirului}
        BEGIN
            new(s.cap);
            s.cap^.ch:= c;
            s.cap^.urm:= nil;
            s.coada:= s.cap [4.2.3.e]
        END
    ELSE
        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 **sir nou destinație** identic cu sursa.
- Un astfel de operator **Copiază**(destinație,sursa) poate fi utilizat cu succes în implementarea operației **ConcatSir** [4.2.3.f].

{Implementarea operatorului ConcatSir utilizând operatorul Copiază}

```

PROCEDURE ConcatSir(VAR u: TipSir; s: TipSir);
VAR temp: TipSir;
BEGIN
    IF u.lungime=0 THEN
        Copiază(u,s)
    ELSE
        IF s.lungime<>0 THEN
            BEGIN
                CreazăSirVid(temp); [4.2.3.f]
                Copiază(temp,s);
                u.coada^.urm:= temp.cap;
                u.coada:= temp.coada;
                u.lungime:= u.lungime+s.lungime
            END
        END
    END
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 înlanțuită**.
- 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 **InsereazăȘir** și **ȘtergeȘir** se realizează în $O(1)$ unități de timp
 - **Căutarea** poziției însă necesită practic $O(n)$ unități.
 - De fapt căutarea **secvențială** a poziției este marele **dezavantaj** al acestei reprezentări care diminuează performanța marii 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}

TYPE TipSir = ARRAY[0..m-1] OF char; [4.3.1.a]

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

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

{Definirea relației de ordonare a două șiruri}

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

- 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 (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}

VAR x,y: TipSir;

[4.3.1.d]

i:= 0;

WHILE (x[i]=y[i])AND(x[i]<>chr(0)) DO i:= i+1; **{O(n)}**

- 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}

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

```

CONST n = {numărul de șiruri din tabela de șiruri};
      m = {numărul maxim de caractere dintr-un șir};
      {sfârșitul unui șir se marchează cu chr(0)}

TYPE TipSir = ARRAY[0..m-1] OF char;
      TipTabela = ARRAY[0..n-1] OF TipSir;      [4.3.1.g]

VAR T: TipTabela;

```

- 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}

```

FUNCTION CautareTabelara(VAR T: TipTabela; VAR x: TipSir;
                        VAR poz: integer): boolean;
VAR i,s,d,mij: integer;
BEGIN

```

```

s:= 0;d:= n;
WHILE s<d DO {căutare binară} [4.3.1.h]
  BEGIN
    mij:= (s+d) div 2; i:= 0;
    WHILE (T[m,i]=x[i])AND(x[i]<>chr(0)) DO i:=i+1;
    IF T[m,i]<x[i] THEN s:= m+1
      ELSE d:= m
    END; {WHILE}
  IF d<n THEN {comparare de siruri}
    BEGIN
      i:= 0;
      WHILE (T[d,i]=x[i])AND(x[i]<>chr(0)) DO
        i:= i+1
      END; {IF}
      poz:= d;
      CautareTabelara:= (d<n)AND(T[d,i]=x[i])
    END; {CautareTabelara}
  
```

- Funcția **CautareTabelară** primește ca parametri de intrare tabela T și sirul 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.
- Se observă că T și x care sunt șiruri, se declară cu **VAR**, pentru a li se transmite (prin stivă) doar adresa.

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}

```

VAR s: ARRAY[0..n-1] OF char;
      p: ARRAY[0..m-1] OF char; [4.3.2.a]
  
```

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

```
typedef unsigned char boolean;
char s[n];
char p[m];
```

/*[4.3.2.a]*/

- 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].

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

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

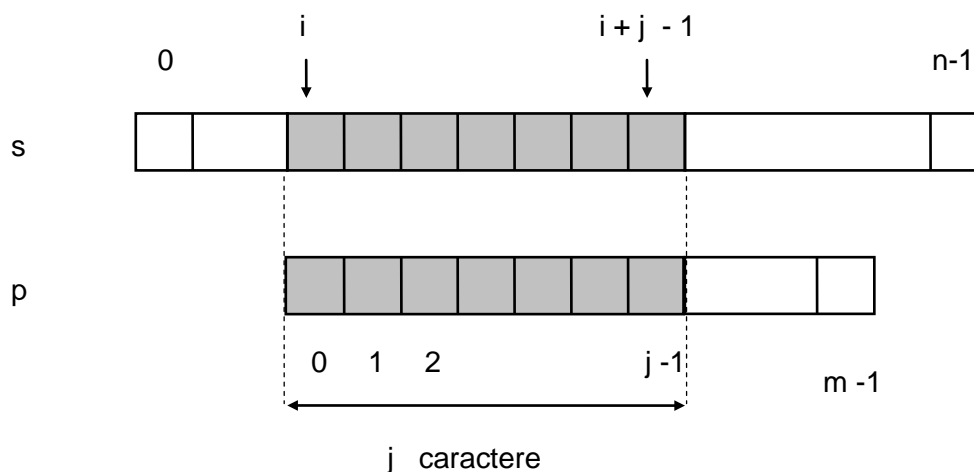


Fig.4.3.2.a. Reprezentarea grafică a predicatului $P(i, j)$

- 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].

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

- 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 - Implementarea predicatului $Q(i)$ }

```
i := -1;
REPEAT
  i := i+1;
  gasit := P(i, m)
UNTIL gasit OR (i=n-m);
```

[4.3.2.d]

- 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 Pascal}

```
FUNCTION CautareDirecta(VAR poz: integer): boolean;
  VAR i, j: integer; {i parcurge caracterele din șir,
                    j parcurge caracterele din model}
  BEGIN
    i := -1;
    REPEAT
      i := i+1; j := 0; {Q(i)}
      WHILE (j<m)AND(s[i+j]=p[j]) DO j := j+1 {P(i, m)}
    UNTIL (j=m)OR(i=n-m);
    poz := i;
    CautareDirecta := j=m;
  END; {CautareDirecta}
```

[4.3.2.e]

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

```
boolean cautare_directa(int* poz)
```

```

{
    int i,j;    /*i parcurge caracterele din șir,
                j parcurge caracterele din model*/
    boolean cautare_directa_result;
    i= 1;
    do {                                              /*[4.3.2.e]*/
        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;
    cautare_directa_result= j==m;
    return cautare_directa_result;
}    /*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

MARGINE

MARGINE

MARGINE

MARGINE

MARGINE

[4 . 3 . 3 . a]

. . .
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;

WHILE (j<m) AND (i<n) DO

[4 . 3 . 3 . b]

BEGIN {Q(i-j) & P(i-j,j)}

WHILE (j>=0) AND (s[i]<>p[j]) DO j:= d;

i:= i+1; j:= j+1

END; {WHILE}

- 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).

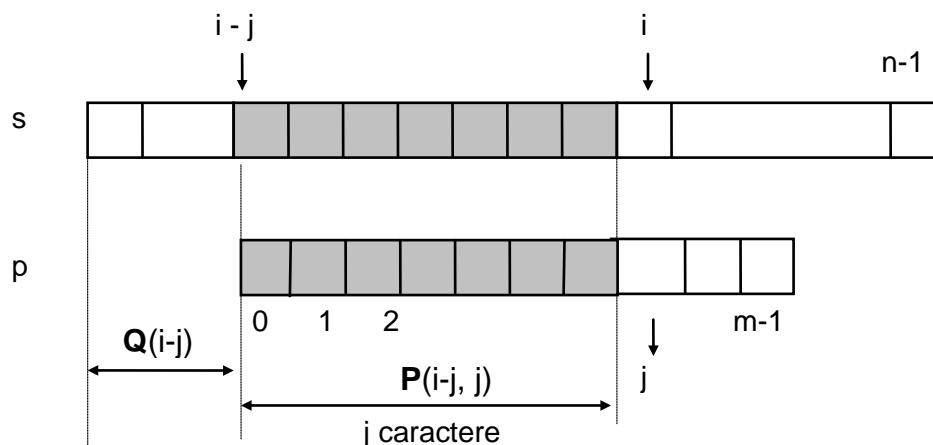


Fig. 4.3.3.a. Reprezentarea predicatelor P și Q în contextul căutării KMP

- 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.

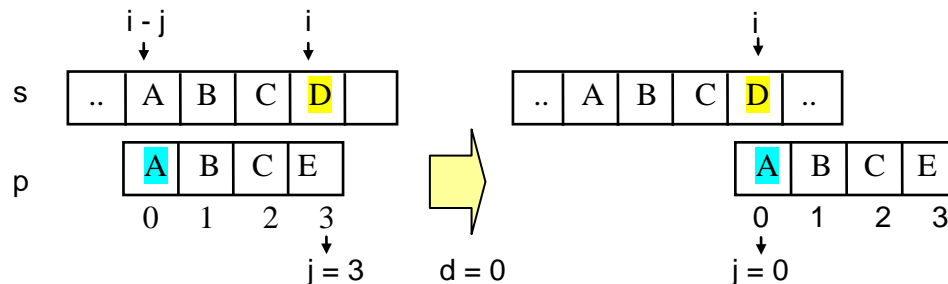


Fig. 4.3.3.b. Deplasarea modelului în cadrul șirului s . Cazul 1.

- 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.

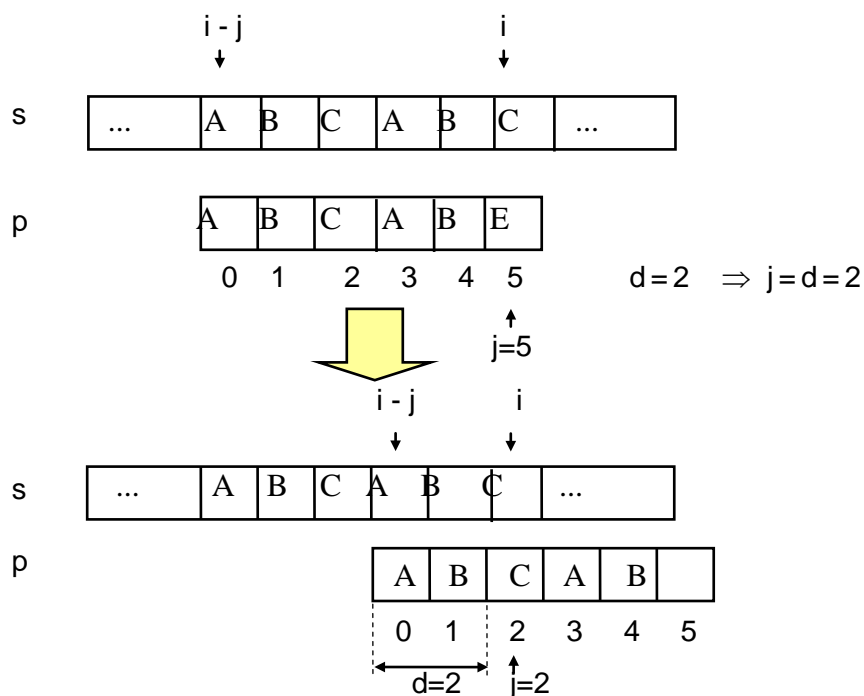


Fig. 4.3.3.c. Deplasarea modelului în cadrul șirului s . Cazul 2.

- 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 cele d poziții ale modelului 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].

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

- 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 valoare a lui 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 propriuzise 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.

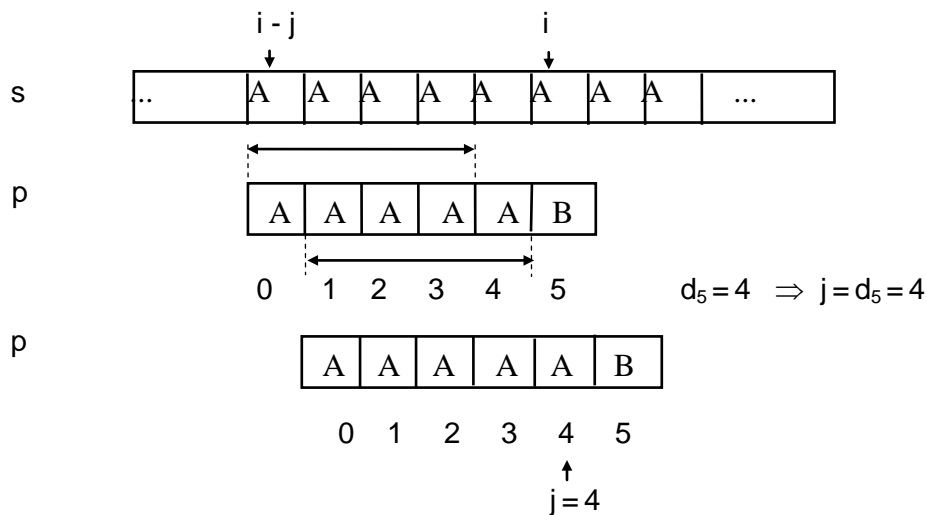


Fig. 4.3.3.d. Determinarea deplasamentului d . Cazul 1.

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

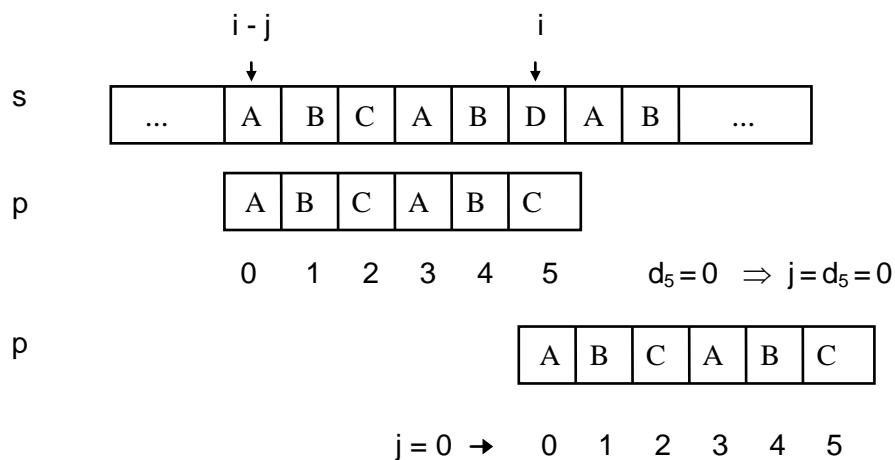


Fig. 4.3.3.e. Determinarea deplasamentului d . Cazul 2.

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

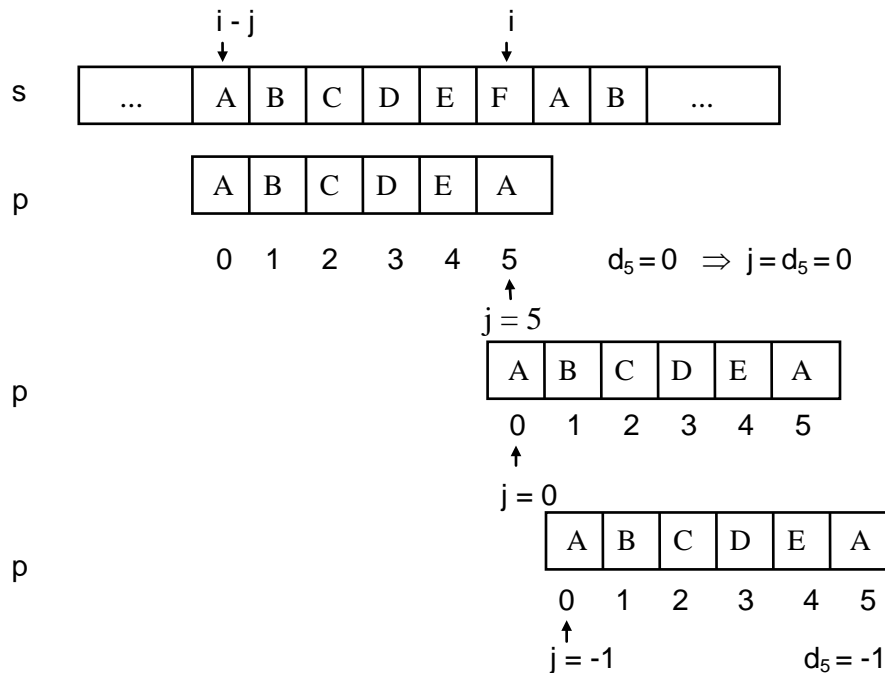


Fig. 4.3.3.f. Determinarea deplasamentului d . Cazul 3.

- 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}

```

CONST mmax={lungime maximă model};
      nmax={lungime maximă șir sursă};

VAR m{lungime model},n{lungime șir}: integer;
    p: ARRAY[0..mmax-1] OF char;{model}
    s: ARRAY[0..nmax-1] OF char;{șir}
    d: ARRAY[0..mmax-1] OF integer;{tabela de deplasări}

FUNCTION CautareKMP(VAR poz: integer): boolean;
    VAR i,j,k: integer;
    BEGIN
[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}
      WHILE j<m-1 DO
        BEGIN
          WHILE (k>=0) AND (p[j]<>p[k]) DO k:= d[k];
          j:= j+1; k:= k+1;
          IF p[j]=p[k] THEN d[j]:= d[k]
            ELSE d[j]:= k
          END;{WHILE}
[4]  i:= 0; j:= 0; {căutare model}
      WHILE (j<m) AND (i<n) DO
        BEGIN

```

[4.3.3.d]

```

        WHILE (j>=0) AND (s[i]<>p[j]) DO j:= d[j];
        j:= j+1;
        i:= i+1;
    END; {WHILE}
    poz:= i-m;
    CautareKMP:= j=m
END;

```

- 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 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 $m \cdot n$.
- De asemenea se subliniază faptul că pointerul i care baleează șirul **nu** merge niciodată înapoi.
 - Spre deosebire 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.
- **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.
- 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.

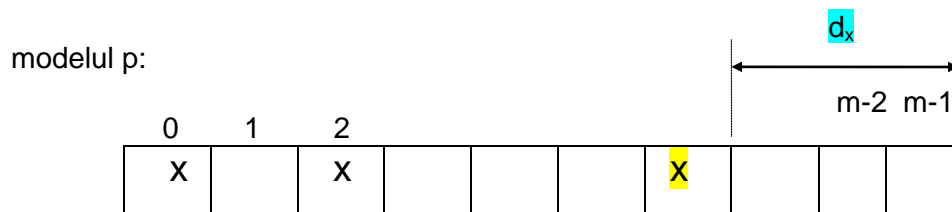


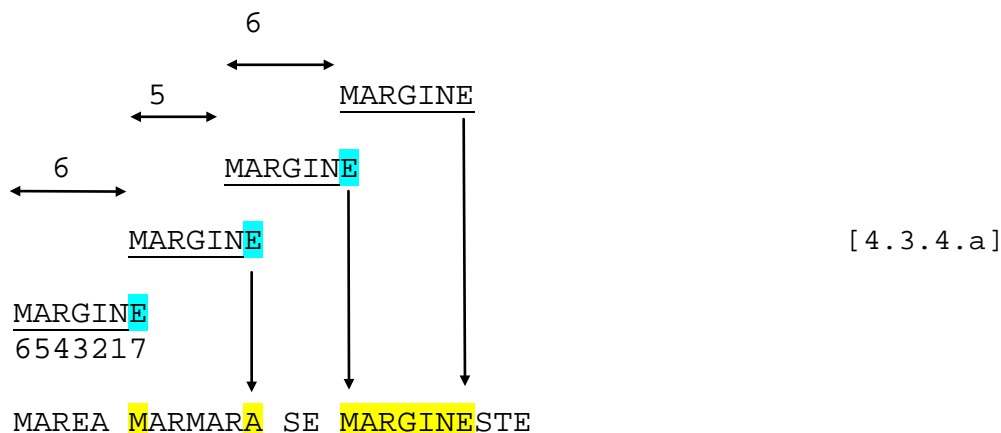
Fig. 4.3.4.a. Determinarea valorii d_x corespunzătoare caracterului x al modelului

- 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.

-
- Diagram illustrating the mapping of a sequence s to a sequence p .
- Sequence s is represented by a horizontal row of cells. The last cell is labeled p_{m-1} . An arrow labeled "nepotrivire" points to the start of sequence s .
- Sequence p is represented by a horizontal row of cells below s . The cells are indexed from 0 to $m-1$, with indices 0, 1, 2, 3, 4, ..., $m-1$ labeled below the cells.
- A double-headed arrow labeled "potrivire" indicates the alignment between the start of sequence p (index 4) and the start of sequence s .

- 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.



- $$\begin{aligned} \mathbf{P}(i, j) &= \mathbf{A}_k: j \leq k < m: s_{i-j+k} = p_k \\ \mathbf{Q}(i) &= \mathbf{A}_k: 0 \leq k < i: \sim \mathbf{P}(k, 0) \end{aligned} \quad [4.3.4.b]$$

- Interpretarea grafică a a noii formulări a predicatelor în cauză apare în figura 4.3.4.c.
- 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).

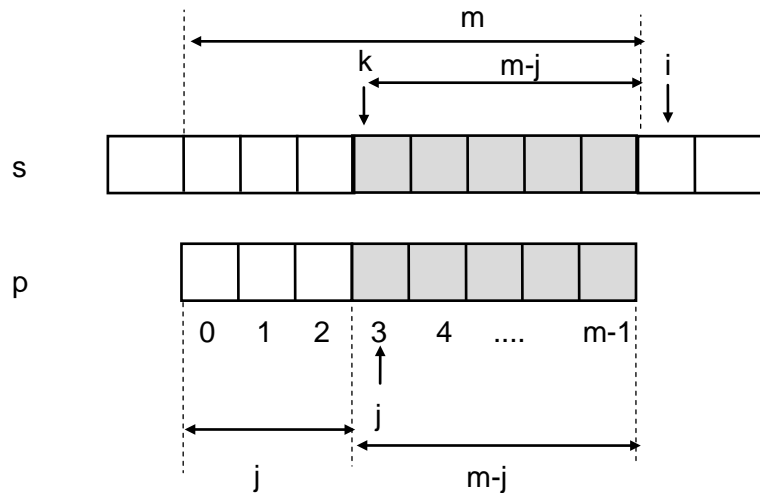


Fig. 4.3.4.c. Interpretarea predicatelor **P** și **Q** în căutarea BM

- Aceste predicate sunt utilizate în formularea următoare a algoritmului BM [4.3.4.c]:

{căutare Boyer-Moore varianta Pascal}

```

i:= m; j:= m;
WHILE (j>0) AND (i<n) DO
  BEGIN {Q(i-m)} [4.3.4.c]
    j:= m; k:= i;
    WHILE (j>0) AND (s[k-1]=p[j-1]) DO
      BEGIN {P(k-j,j)&(k-j=i-m)}
        k:= k-1;
        j:= j-1
      END;
    IF j>0 THEN i:= i+d[ORD(s[i-1])]
  END; {căutare Boyer-Moore}

```

/*căutare Boyer-Moore varianta C*/

```

i= m; j= m;
while ((j>0) && (i<n))
{
  /*Q(i-m)*/ /*[4.3.4.c]*/
  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;
    }
    if (j>0) i= i+d[s[i-1]];
}
/*cautare Boyer-Moore*/
/*-----*/

```

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

```

CONST mmax={lungime maximă model};
      nmax={lungime maximă șir sursa};           [4.3.4.c]

VAR m{lungime model},n{lungime șir}: integer;
    p: ARRAY[0..mmax-1] OF char;{model}
    s: ARRAY[0..nmax-1] OF char;{șir}
    d: ARRAY[0..255] OF integer;{tabela de deplasări}

FUNCTION CautareBM(VAR poz: integer): boolean;
VAR i,j,k: integer;
BEGIN
    *citire șir {n este lungimea curentă a șirului}
    *citire model {m este lungimea curentă a modelului}
    FOR i:=0 TO 255 DO d[CHR(i)]:= m;{inițializare tabelă}
    FOR j:=0 TO m-2 DO d[p[j]]:= m-j-1;{precompilare model}
    i:= m; j:= m; {căutare model}
    WHILE (j>0) AND (i<=n) DO
        BEGIN
            j:= m;
            k:= i;
            WHILE (j>0) AND (s[k-1]=p[j-1]) DO
                BEGIN
                    k:= k-1;
                    j:= j-1;
                END;
            IF j>0 THEN i:= i+d[ord(s[i-1])]
        END; {WHILE}
    poz:= i-m; {poz:=k}
    CautareBM:= j=0

```

END; {Cautare BM}

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

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

enum { mmax = 100 /*lungime maximă model*/,
       nmax = 200} /*lungime maximă șir sursă*/;
int m/*lungime model*/,n/*lungime șir*/;
char p[mmax];          /*model*/
char s[nmax];          /*șir*/
int d[256];            /*tabela de deplasări*/

boolean cautare_bm(int* poz)
{
    int i,j,k;          /*[4.3.4.c]*/

    /*citire sir;      n este lungimea curentă a șirului
    /*citire model;    m este lungimea curentă a modelului
    boolean cautare_bm_result;
    for(i=0; i <=255; i++) d[i]= m;/*inițializare tabelă*/
    for(j=0; j<=m-2; j++) d[p[j]]= m-j-1;/*precompilare
                                                model*/

    i=m; j=m; /*căutare 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]];
    } /*WHILE*/
    *poz=i-m; /*poz=k*/
    Cautare_bm_result= j==0;
    return cautare_bm_result;
} /*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 șir 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.

4.4. Aplicații propuse

Aplicația 4.4.1

Se cere să se implementeze **TDA șir** în cele trei modalități prezentate în cadrul capitolului de față. Pentru cele trei implementări, se cere să se realizeze o analiză comparată a performanțelor fiecărui operator în termenii operatorului O.

Aplicația 4.4.2

Se cere să se realizeze un studiu al performanțelor căutării tabelare care să evidențieze următoarele aspecte:

- Influența metodei de căutare propriu-zise (liniară sau patrică) în raport cu dimensiunea tabelii.
- Influența lungimii cheii de căutare asupra performanței căutării.

Rezultatul studiului va fi construcția câte unui profil al algoritmului de căutare în fiecare din circumstanțele impuse. Se vor comenta comparativ rezultatele obținute.

Aplicația 4.4.3

Se cere să se redacteze un program interactiv destinat studiului performanțelor metodelor de căutare în șiruri care implementează următoarele comenzi:

S – introducere șir sursă. Introducerea se poate face direct prin tastare, prin citirea unui fișier text sau prin generarea aleatoare a unui șir de caractere cu lungimea precizată care conține pe ultima poziție un text dorit.

P - introducere șir model.

D - căutare model prin metoda directă.

K - căutare model prin metoda Knuth-Morris-Pratt.

B - căutare model prin metoda Boyer-Moore.

T - terminare.

Pentru fiecare dintre metodele de căutare, la terminare se afișează:

- Poziția coincidenței în cadrul șirului sursă.
- Numărul de comparații de caractere efectuat.
- Timpul necesar realizării căutării.