

Note asupra realizării unui algoritm paralel de căutare a unui șablon într-un text

Lect. dr. Teodor Florin Fortiș
Universitatea de Vest, Timișoara

ABSTRACT. The purpose of this article is to realize an analysis for two methods that can be considered in order to create a parallel algorithm for pattern searching into a text. First, we are dealing with "Divide-and-Conquer" type approaches. Next, we are considering another type of approach, based on the decomposition of the initial problem into several simpler problems and with different initial data, but similar solutions. An algorithm for the later situation is sketched, in order to reveal some of the (possible) features of algorithms of this type.

Key words: pattern matching, string searching, text algorithms.

1 Probleme clasice de căutare

Problemele de căutare a șirurilor de caractere se dovedesc a fi un subiect extrem de important într-o varietate de domenii, în principal legate de procesarea textelor. Multe sisteme software se bazează pe funcționarea unor algoritmi de căutare eficienți, aceștia oferind de cele mai multe ori nu numai un sprijin sistemelor software dar chiar și modele de programare.

Informația de bază este depozitată, în multe domenii ale științelor, într-o manieră lineară, similară șirurilor de caractere. Acest mod de depozitare a informației este evident în domenii filologice, dar și în alte domenii incluzând informatica (unde informația poate fi stocată în fișiere de dimensiuni mari) sau biologia (unde moleculele pot fi descrise oferind o secvență de nucleotide și amino-acizi). Cantitatea mare de date care urmează să fie procesată presupune existența unor metode eficiente de identificare a informației.

O problemă de căutare presupune identificarea primei/unei apariții (sau a tuturor aparițiilor) a unui șir de simboluri dintr-un alfabet, numit șablon, într-un text. Atât șablonul cât și textul sunt construite peste o mulțime finită de caractere, numită alfabet.

În funcție de modul în care sunt oferite informațiile de bază, pot fi realizate mai multe abordări ale acestei probleme: cunoașterea prealabilă a șablonului face posibilă utilizarea unui automat sau a anumitor proprietăți ale șablonului, pe când cunoașterea prealabilă a textului face posibilă utilizarea unor proprietăți statistice ale acestuia în scopul îmbunătățirii performanțelor algoritmilor.

Funcționarea algoritmilor bazați pe cunoașterea prealabilă a șirului de căutare poate fi aproximată astfel:

- textul este parcurs prin intermediul unei ferestre, a cărei dimensiune este egală cu dimensiunea șablonului;
- fereastra este aliniată inițial la stânga cu textul, realizându-se o comparație a conținutului ferestrei cu conținutul șablonului. Această operație este numită încercare ;
- după fiecare încercare este realizată o translatare a ferestrei spre dreapta, algoritmul pregătindu-se pentru o nouă încercare;
- fereastra este mutată către dreapta până la atingerea marginii drepte a textului în care se realizează căutarea.

Aplicarea directă a acestei metode oferă algoritmul forței brute, capabil să identifice toate aparițiile unui șablon într-un text într-un timp $O(m \times n)$, unde m este mărimea șablonului iar n este dimensiunea textului. Creșterea vitezei algoritmului este posibilă prin realizarea unor îmbunătățiri aduse felului în care este realizată comparația în interiorul ferestrei: de la stânga la dreapta, de la dreapta la stânga, într-o ordine bine precizată sau fără restricții.

În prima categorie intră algoritmii care se bazează pe funcții de hashing (Harrison, Karp-Rabin) sau recunoașterea unui șir de caractere cu ajutorul unui automat (Morris-Pratt, Knuth-Morris-Pratt). A doua categorie conține unul dintre cei mai eficienți algoritmi utilizați în practică (Boyer-Moore) sau variante ale acestuia (Colussi, Turbo-BM, și altele).

A treia categorie de algoritmi (în care comparația este realizată într-o ordine bine precizată) oferă algoritmii cu cele mai bune performanțe teoretice, (Galil-Giancarlo, Galil-Seiferas, Crochemore-Perrin, etc.). Aceștia se bazează de cele mai multe ori pe algoritmi din primele două categorii, oferind o modificare simplă a acestora.

2 Abordări paralele pentru căutarea unui șablon într-un text

2.1 "Divide-and-conquer"

Cea mai simplă și, probabil, cea mai directă modalitate de construire a unui algoritm paralel de căutare a unui șablon într-un text presupune utilizarea conceptelor "Divide-and-conquer" pentru construirea subproblemelor care urmează să fie abordate în paralel. În acest fel, problema de căutare principală este descompusă în mai multe probleme de căutare pe bucăți distincte ale textului inițial. Fiecare dintre procesele participante la rezolvarea paralelă va avea de realizat căutarea aceluiași șablon în bucăți diferite de text.

În mod ideal, în condițiile în care problema principală este descompusă în n subprobleme, viteza de căutare ar putea crește de cel mult n ori. Cu toate acestea, nici unul dintre procesele participante nu va putea lua nici un fel de decizie la capetele textului pe care-și desfășoară activitatea, astfel încât încercarea de a păstra această viteză superioară va fi, inevitabil, însoțită de posibilitatea de a pierde una sau mai multe apariții ale șablonului.

Pentru a trece peste această dificultate pot fi avute în vedere următoarele soluții:

1. Extinderea uniformă a textului fiecărui proces cu câte $\left\lceil \frac{m}{2} \right\rceil$ simboluri, unde m este lungimea șablonului, în ambele capete. În felul acesta, textele pasate către două procese vecine se vor putea suprapune pe o lungime de cel mult m simboluri.
2. Realizarea de căutări parțiale (ipoteze) la capetele șirului, urmând ca acestea să fie confirmate în momentul în care căutarea în zona vecină este încheiată.

Alegerea metodei de căutare poate fi realizată în funcție de modul în care are loc transferul de date între procesul principal și cele n procese desemnate pentru rezolvarea subproblemelor.

În situația în care procesele au acces la o memorie comună prin intermediul căreia pot accesa textul în care se realizează căutarea, o alegere potrivită pentru rezolvarea fiecăreia dintre subproblemele de căutare este unul dintre algoritmii cu viteză de căutare mare, cum ar fi Boyer-Moore, Colussi, Turbo BM.

Dacă, însă, este necesară transmiterea integrală a datelor către procesele participante, procesul central va putea realiza o estimare a frecvențelor caracterelor în timpul construirii pachetelor de date care urmează să fie transmise.

În acest mod ar putea fi ales un algoritm în care comparația caracterelor din șablon cu caracterele din fereastră este realizată într-o ordine bine precizată, alta decât de la stânga la dreapta sau de la dreapta la stânga.

2.2 Transformarea condițiilor problemei

Abordarea avută în vedere în acest articol oferă un punct de vedere ușor diferit. Paralelismul este introdus în acest caz prin transformarea uniformă atât a spațiului de lucru (textului) cât și a șablonului utilizat. În această situație problema de căutare este redusă la mai multe subprobleme de căutare diferite, folosind șabloane și texte modificate prin aplicarea unei proiecții.

Identificarea aparițiilor unui șablon într-un text este realizată prin combinarea rezultatelor întoarse de fiecare dintre subprobleme, această operație fiind realizată îndată ce este semnalată o posibilă potrivire a șablonului în text. Totodată, în acest mod este posibilă accelerarea vitezei de căutare, diferitele procese participante la proces având astfel capacitatea de a "sări" peste bucăți mari de text, dacă este posibil.

Transformarea textului și a șablonului este realizată într-o manieră extrem de simplă: fie V alfabetul folosit pentru scrierea textului, $V_1 \subset V$ o submulțime a lui V iar \bar{b} un simbol (numit blank) care nu se găsește în alfabetul V . Definim aplicația $pr_{V_1} : V \rightarrow V_1 \cup \{\bar{b}\}$ prin

$$pr_{V_1}(a) = \begin{cases} a, & a \in V_1 \\ \bar{b}, & a \notin V_1 \end{cases}$$

Pentru început, descompunem alfabetul V utilizat pentru descrierea textului într-un set de k alfabete disjuncte, formate, de preferat, din simboluri cu frecvențe cumulate de apariție apropiate. În plus, toate simbolurile utilizate în aceste mulțimi trebuie să apară în șablonul principal. Această descompunere limitează numărul de alfabete disjuncte (și, implicit, numărul de procese participante la rezolvarea problemei) la numărul total de simboluri utilizate pentru scrierea șablonului.

Fiecare dintre procesele P_i cu $i = 1, 2, \dots, k$, participante la algoritmul de căutare recepționează proiecțiile peste alfabetul V_i ale textului și ale șablonului utilizat. Mai departe, activitatea fiecărui proces P_i presupune căutarea proiecției șablonului original în proiecția textului original, fiecare apariție fiind numită ipoteză. Fiecare proces se orientează așadar spre emiterea de ipoteze și verificarea ipotezelor emise de alte procese.

2.2.1 Caracteristicile metodei

Subliniem, în cele ce urmează, câteva dintre caracteristicile care pot fi deduse din această metodă și din algoritmul (vezi Secțiunea 2.2.2) utilizat pentru funcționarea fiecăruia dintre procesele P_i cu $i = 1, 2, \dots, k$.

Viteza de căutare

Procesele care realizează o căutare rapidă într-o porțiune a textului propriu pot influența accelerarea funcționării celorlalte procese în mai multe moduri:

- după emiterea unei ipoteze de către un astfel de proces, celelalte procese sunt "obligate" să avanseze la poziția ipotezei pentru verificarea acesteia;
- la verificarea unei ipoteze emisă de către un alt proces, P_j , acesta va putea avansa la poziția curentă a unui proces mai rapid.

Prin urmare procesele mai lente pe anumite porțiuni sunt obligate să realizeze salturi în funcționare pentru a ajunge la poziția proceselor mai rapide.

Transferul eficient al datelor

Fiecare proces P_i cu $i = 1, 2, \dots, k$, participant la algoritmul de căutare își desfășoară activitatea în condiții deosebite: majoritatea simbolurilor din textul și din șablonul utilizat (proiecția pr_{V_i} a textului și a șablonului original) sunt simboluri blank, \bar{b} . Textul poate fi privit ca o matrice rară, cu numeroase goluri, posibil continue, marcate prin simboluri blank. Având în vedere această calitate a datelor, poate fi implementată o metodă pentru

optimizarea modului în care urmează să fie transferate datele problemei către procese.

Adaptabilitate

Pregătirea fiecărei proiecții se realizează prin intermediul unui proces specializat. Prin intermediul acestui proces este posibilă culegerea informațiilor legate de frecvențele de apariție ale caracterelor. În acest mod, în situația în care datele nu sunt transmise către procese într-un singur transfer, este posibilă modificarea celor k mulțimi și a proiecțiilor asociate astfel încât distribuția simbolurilor din cele k mulțimi să fie cât mai uniformă posibil.

Funcționare sincronă vs. funcționare asincronă

Implicit, algoritmul presupune un comportament sincron al proceselor P_i cu $i = 1, 2, \dots, k$. Totuși, în ipoteza că este folosit un model client-server pentru a suporta verificarea ipotezelor, este posibilă bufferizarea ipotezelor emise de cele k procese, lăsând problema verificării ipotezelor și a accelerării vitezei proceselor participante în seama unui proces central (serverul, în această ipoteză). În acest caz comportamentul celor k procese, P_i cu $i = 1, 2, \dots, k$, este asincron.

2.2.2 Descrierea algoritmului

Pentru descrierea algoritmului, vom presupune că fiecare dintre cele k procese participante, P_i cu $i = 1, 2, \dots, k$, utilizează o metodă M_i pentru identificarea următoarei apariții a șablonului în text.

Comunicarea între două procese P_a și P_b presupune transmiterea următoarelor tipuri de mesaje:

- verificare ipoteză: este transmisă poziția ipotezei;
- confirmare ipoteză: este returnată poziția ipotezei;
- infirmare ipoteză: este returnată cea mai avansată poziție dintre poziția curentă și poziția următoare poziției ipotezei.

Fiecare dintre cele k procese participante, P_i cu $i = 1, 2, \dots, k$ își desfășoară activitatea pe textul și pe șablonul propriu în felul următor:

- 1) Procesul P_i utilizează metoda M_i pentru a identifica următoarea apariție a proiecției proprii a șablonului în proiecția textului;
- 2) Dacă este identificată o apariție a șablonului în text, atunci
 - a) Dacă procesul P_i este în timpul verificării unei ipoteze și poziția ipotezei verificate este identică cu poziția ipotezei obținute, emite mesajul de confirmare ipoteză;
 - b) În caz contrar,
 - i) emite câte o cerere de verificare a ipotezei către fiecare proces P_j , cu $j = 1, 2, \dots, k, j \neq i$.
 - ii) Așteaptă confirmarea ipotezei de la fiecare proces P_j , cu $j = 1, 2, \dots, k, j \neq i$.
 - iii) Dacă ipoteza este confirmată, comunică apariția procesului central, apoi trece la Pasul 1.
 - iv) Dacă ipoteza este infirmată de un proces P_j , trece la poziția returnată în mesajul de infirmare, apoi trece la Pasul 1.
- 3) Dacă este recepționată o cerere de verificare a unei ipoteze din partea unui proces P_j , atunci:
 - a) Dacă poziția curentă depășește poziția ipotezei, întoarce poziția curentă către procesul P_j în mesajul de infirmare a ipotezei.
 - b) Dacă poziția ipotezei depășește poziția curentă, avansează pe poziția ipotezei, apoi treci la Pasul 1.
- 4) Dacă este epuizat spațiul de lucru al procesului curent, emite o cerere fictivă de verificare a unei ipoteze pentru poziția finală, pentru a determina terminarea tuturor celorlalte procese. Încheie execuția procesului P_i .

Bibliografie

- [AG92] **Apostolico, A., Giancarlo, R.**, *The Boyer-Moore-Galil string searching strategies revisited*, SIAM Journal on Computing, **15** (1), 98--105, 1992.

- [Aho90] **Aho, A.V.**, *Algorithms for Finding Patterns in Strings*, Handbook of Theoretical Computer Science, Volume A, Algorithms and complexity, J. van Leeuwen ed., Cap. 5, pp **255 - 300**, Elsevier, Amsterdam, 1990.
- [BM77] **Boyer, R.S., Moore, J.S.**, *A fast string searching algorithm*, Communications of the ACM, **20** (1977), 762-772.
- [CR94] **Crochemore, M., Rytter, W.**, *Text Algorithms*, Oxford University Press, 1994.
- [KMP77] **Knuth, D.E., Morris, Jr, J.H., Pratt, V.R.**, *Fast pattern matching in strings*, SIAM Journal on Computing, **6** (1), 323---350, 1977.
- [Ste94] **Stephen, G.A.**, *String Searching Algorithms*, World Scientific, 1994.