

## **ALGORITMI ZA ISKANJE REZERVIRANIH BESED**

M. GAMS

UDK: 519.688

## **INSTITUT „JOŽEF STEFAN“, JAMOVA 39, LJUBLJANA**

U dlaniku so opisani algoritmi za iskanje rezerviranih besed. Glavna algoritma sta binarno iskanje in iskanje z rezerviranimi tabelami. V dlanku so opisane izboljšave teh dveh osnovnih algoritmov. Podana je tudi njihova razvina in prostorska analiza.

**ALGORITHMS FOR SEARCHING FOR RESERVED WORDS.** This article describes algorithms for searching for reserved words. Basic algorithms are binary search and hash search. This article describes basic algorithms with their improvements and provides time and memory requirements analysis.

### **1. Used**

U računalniških programih posostaje rezervirane besede. To so besede, ki imajo potreben rezerviran pomen. Take besede rezervamo v prevajalnikih, urejevalnikih, operacijskih sistemih itd. Rezervirane besede so del teksa, npr. programa v programske jeziku, zato morajo radunalniški programi vsebovati postopek za iskanje rezerviranih besed. Osnovna algoritma za iskanje rezerviranih besed sta binarno iskanje in iskanje z razrednjanimi tabelami. Z izboljšavami teh dveh algoritmov je možno dosegati prenosljive rezultante.

## 2. Open algoritmos

Najprej bomo opisali osnovne algoritme in nato njuno izboljšave. Algoritme bomo opisali v Pascalju. Kjer pa bomo prekoračili izrazna možnosti tega jezika, so ta rezaberljivo označeno. V tem poglavju bomo govorili tudi o konvergencijsnosti algoritmov. S tem pojmom bomo opredelili povrtevno število primerjavi neznane besede z rezerviranimi.

#### **3.4. Síndrome de Asperger**

### Open algorithms

```

lower := word[higher] := bottom
found := false
while (lower <= higher) and not found do
begin
  middle := (lower + higher) div 21
  if words[middle] > unknownWord
    then higher := middle
  else if words[middle] < unknownWord
    then lower := middle
  else found := true
end

```

	WORDS
lower =>	LAND ARRAY ...
middle =>	...
higher =>	INETH

Operasi

- na vsebini obvezuječih Pascalogovih prevajalnikov operatorja "<" in ">" nista implementirana za prirejanje sestavljenih timov.
  - rezervirane besede v tabeli "WORDS" morajo biti urejene po abecednem redu.
  - sestavljalniki "lower" in "higher" morata biti pred zabeščenim izvajanjem binarnega iskanja prirejeni na vrh, oziroma dno tabele.

Konkretnost algoritma za binarne iskanja je reda velikosti los(n), kjer je n število rezerviranih besed. Ker je rezerviranih besed običajno nekaj deset, je  $los(n) < 5$  ali 6.

Najbolj običajan trik za potrditev iskanja je grupiranje rezerviranih besed v skupine. Eden principov za grupiranje v skupine je po dolžini besed, tako da se znotraj iste skupine besede imajo dolžine urejene po abecednem redu. Algoritem za binarno iskanje je potreben samo malenkostno spremeniti. Binarno iskanje vrstimo v tabeli "WORDB" takor da

PRED IZVAJANJEM OSNOVNega ALGORITMA PRIREDIMO SPREMEMljivKama "lower" in "higher" VREDNOSTI Iz VNAPREj ZENERIRANE TABELE, ki VSEBUJE INDEKS ZAETKA in konca skupin besed iste dolzine. Na Primer, če je "1" dolzina neznane besede, bo "lower := lensht[1]" in "higher := lensht[lj] - 1". Primer za JEZIK PASCAL JE SNETKIChno PRIKAZAN Malo naprej.

LENGTH	WORDS
2	I00
3	IIF
4	IN
5	IGF
6	IOR
7	IIO
8	IAND
9	...
10	...

Za realne PRIMERe se KOMPLEKSNOST ALGORITMA Priblizno dvakrat zmanjša in je med 2 in 3.

Druš nacin GRUPIRANJA je PO PRVI ali PO PRVIH dveh CRKAH. V prvem PRIMERU je tabela za določanje indeksov v tabeli WORDS razdeljena na 26 in v drugem na 676 delov. Dodatno potrebimo 26 ali 676 celio seominar, iskanje neznane besede pa ima Kompleksnost med 2 in 1. Možne so še razne kombinacije, npr. grupiranje po dolžini in po prvi CRKI itd.

TWOCHARS	WORDS
AA	I00
AB	I01
AC	I01
..	...
..	...
AN	I11
AO	I01
AP	I01
AQ	I01
AR	I21
..	...
..	...

### 2.2. ISKANJE Z RAZPRESENIMI TABELAMI

OSNOVNI ALGORITEM ISKANJA Z RAZPRESENIMI TABELAMI TEMELJI NA FUNKCIJI "h", ki preslikava vse besede v indeks tabel. Kompleksnost takih algoritmov je običajno kar konstanta. Kompleksnost je 1 kadar je funkcija injektivna ( $h(w1) \neq h(w2)$ ). Običajno imajo razperšitvene funkcije lastnosti, ki so v povprečju blizu injektivnosti, tako da se lahko prioritizira kompleksnost 1. Druga pomembna lastnost razperšitvenih funkcij je to, da običajno razperšijo podatke v precej vedno tabelo, oziroma da je "tabeli precej praznih mest. Od tod tudi ime "razperšene tabele".

Kadar je razperšitvena funkcija injektivna, rečemo, da je "popolna" (perfect). Kadar je razperšitvena funkcija surjektivna, rečemo, da je "minimalna". Minimalna razperšitvena funkcija ima v našem PRIMERU tabelo "WORDS" polno zasedeno.

V zadnjih letih je nekaj odkritij povečalo zanimanje za minimalne ali skoraj minimalne popolne razperšitvene funkcije [1,2,3]. V [1] je opisana minimalna razperšitvena funkcija implementirana na rezerviranih besedah jezika Pascal. Razperšitvena funkcija je oblike

```
h(Word) = value[Word[1]] + lensht + value[Word[lensht]]
```

"lensht" je dolžina neznane besede "Word" "Word[1]" je prva CRKA besede "Word" "value[letter]" je izbrana vrednost CRKE "letter"

Obitno funkcija ne more biti popolna. Kadar obstajata dve rezervirani besedi enake dolžini, ki imata paroma enaki prvi in zadnji CRKI. V [2] je pokazano, da obstaja več drugačnih protiperimerov, torej postopek iskanja razperšitvenih funkcij na zornji način ne jamči rešitve. V [3] je omemljena metoda, ki nima te slabosti, vendar je tako zenerirana funkcija časovno zahtevnejša.

PROBLEM ISKANJA RAZPERŠITVENE FUNKCIJE JE PROBLEM PRIREDJANJA VREDNOSTI POSAMEZNIM CRKAM. Denimo, da imamo samo 20 različnih začetnih in končnih CRK rezerviranih besed, ter da vsaki lancu priredimo največ 20 različnih vrednosti. Kompleksnost preprostega algoritma, ki preračuna vse variente, je  $20^{*}20$  (dvajset na dvajset), torej praktično neizvedljivo. Napišimo ga:

```
repeat
    if overlap then letter := pred(letter)
    else besin letter := succ(letter);
        value[letter] := -1
    end;
repeat
    value[letter] := value[letter] + 1;
    checkOverlaps(overlap,value,words)
until not overlap or
    (value[letter]>maxValue)
until not overlap and (letter = 'z')
```

OSNOVNI ALGORITEM LAKO POKITIMO Z RAZNIMI TRIKI. NAMESTO ZAPOREDNESA IZBIRANJA CRK ZAČNEMO Z NAJBOLJ POGOSTIMI. Nadaljnje izboljšave niso objavljene v literaturi. Besrajo so v PRAKSI NAJBOLJ UPORABLJANE. NAJBOLJ OBINKOVITO POKITITEV IMENOVANA "superbacktrack". Kadar se dve rezervirani besedi prekrivata tako, da ju ne moremo ložiti s spremembo vrednosti tekotih CRK, moramo spremeniti vrednost CRKE, ki nastopa v priredjanju in smo ji zadnji doležili vrednost. Posledno si PRIMER:

```
value[letter]
-----
e : 0
r : 1
t : 4
f : ?    ;;; TEKOČA CRKA
w1 = For
w2 = File
h(w1) = value[f] + 3 + value[r] + value[e] + 4
h(w2) = value[f] + 4 + value[e] + value[r] + 4
```

"Superbacktrack" v takem PRIMERU NOVEČA VREDNOST CRKE "r" za 1.

AUTOR JE OPAZIL, ŽE EN KORISTEN RECEPT. KO ZASLEDUJEMO IZVAJANJE ALGORITMA, OPAZIMO, DA SO NEKATERA ZAPOREDNA CRKA ZELO NEPRIJETNA. TAKRAT PROGRAM SPREMENI VREDNOST CRK KORAK ZA KORAKOM BREZ "superbacktracka". TAKRAT JE NAJBOLJ PREKHINJITI IZVAJANJE IN SPREMENITI VREDNOST PRIREDJANJA VREDNOSTI CRKAM.

Tako je za generiranje minimalne popolne razprtivene funkcije za Pascalove rezervirane besede avtorjev program potrebil le nekaj sekund. Posledno si primer, ko vendarjej izberemo najmanjšo možno vrednost razprtivene funkcije: Je 3:

VALUE	WORDS
1 IEI 01	I
2 IPI 01	I
3 IDI 01	IEND
4 ITI 11	ELSE
5 INI 31	TYPE
6 IDI 81	PACKED
7 IRI 81	INOT
8 IFI 111	ITHEN
9 ILI 111	IPROCEDURE
10 IA 1151	DO
11 IC 1191	TO
12 II 1191	RECORD
13 IM 1231	REPEAT
14 IV 1251	DOWNTO
15 IW 1141	FILE
16 IB 1211	OR
17 IG 1191	INIT
18 IH 1151	AND
19 IS 1311	WHILE
20 IU 1201	IFOR
21 IY 1171	IOF
22 IJ 1 01	IFUNCTION
23 IK 1 01	ICASE
24 IO 1 01	IIN
25 IX 1 01	ICONST
26 IZ 1 01	IMOD
27 -----	ILABEL
28 -----	IDIV
29 -----	IBEGIN
30 -----	IPROGRAM
31 -----	IGOTO
32 -----	IIF
33 -----	IWITH
34 -----	IVAR
35 -----	ISET
36 -----	IUNTIL
37 -----	IARRAY

Za Pascalove rezervirane besede znano enostavno in hitro določiti minimalno popolno razprtivo funkcijo. Kaj pa če imamo v rezerviranih besedah tudi besede kot "odr" in "car" ali "fi" in "if"? Obično moramo izbrati drugačno obliko razprtivene funkcije. Nato poskusimo generirati tabelo "value". Ker nam postopek ne jandi rešitev, prav lahko poskušamo zmanjšati. V takih primerih je smiselno najprej iskati popolno razprtiveno funkcijo, ki ni minimalna. Tabela "WORDS" je zdaj dimenzija  $N + NEKAJ$ , kjer je  $N$  število rezerviranih besed. Ko najdemo eno rešitev, lahko zmanjšujemo "NEKAJ", dokler ne najdemo rešitev. Da ne bi trošili preveč spomina, si pomagamo že z naslednjim trikom:

VALUEH	WORDS
MIN	I 11 → ...
MIN + 1	I 01
MIN + 2	I 21
...	...   ...
MIN+NEKAJ	I 351

Tako potrebimo se  $N + NEKAJ$  celic spomina.

### 3. Podrobnejša besovna in prostorska analiza algoritmov

Dot sedaj smo sicerili le o povprečnem številu potrebnih primerjanj neznane besede z rezerviranimi. To smo označili s "kompleksnostjo algoritma". Sledi podrobnejša analiza algoritmov. Za eno operacijo bomo šteli eno primerjanje, eno prirejanje, eno računska operacija ( $+, \cdot, \dots$ ), en dosesa polja itd. Te operacije besovno niso enakovredne, zato so ocene okvirne in odvisne od računalnika. Priuzeli bomo verda en znak zavzema eno besedo spominata da je rezervirana beseda sestavljena iz 10 znakov in da je vseh rezerviranih besed 32.

Ocenimo število osnovnih operacij za primerjanje besed med seboj:

```
while (a[i] == b[i]) and (i < dimWord) do
    i := i + 1;
```

Imamo:

```
1 prirejanje (i=)
1 rač.op. (+)
1 los.op. (and)
2 primerjanje (=,<)
2 dosesa polja (i[i])
+1 za kontrolni konstrukt "while"
```

8 osnovnih operacij

Za enaki besedi nadaloma porabimo toliko primerjanj kot je dimenzija besede (maksimalna dovoljena dolžina). Za različne besedi porabimo bistveno manj primerjanj, v povprečju bomo ta naša ocena uvelj 1,5.

A - neznana beseda je rezervirana  
B - neznana beseda ni rezervirana

ALGORITEM	#TEVILLO OPERACIJ	CELIC SPOMINA
I binarne	180	125
I + po dolžini	140	80
I + po prvi	120	55
I srki	20	20
I + po prvi + dolž.	55	10
I dveh srkah	55	10
I minim.por.	55	12
I razrez. f.	57	14
I popolna	57	14
I razrez. f.	57	14

Opombe:

- ocene so okvirne, ker so zelo odvisne od oblike rezerviranih besed itd.

- Konkretno izvedbe algoritmov lahko malo pohitrimo z enostavnimi triki, ki so v besovni analizi upoštevani. Tpr. kadar v povprečju dostikrat prirejamo enake besede, bomo primerjanje naredili tako:

```
while (unknownWord[i] == WORDS[indexs[i]]) and
    (WORDS[indexs[i]] <> ' ') do
    i := i + 1;
```

Tako primerjanja precej pohitrimo.

- v oceni je bila privzeta taka oblika razpršitvene funkcije kot za Pascalove rezervirane besede.

- v praktični uporabi algoritmov, npr. pri prevajalnikih, odraže večji del potrošenega časa na branje, tako da so prihranki manjši.

#### Praktični rezultati:

Skupina študentov (M.Bradeško, B.Oremuž, M.Hrvatin, T.Kokalj) je merila razmerje med osnovnim binarnim iskanjem in pomolno razpršitveno funkcijo za Pascalove rezervirane besede. Razmerje se je ukalo od 1:8 do 1:9 sledje na način generiranja testnih besed od samih rezerviranih do samih nerezerviranih. To razmerje ustreza teoretično izračunanemu.

Rad bi se zahvalil tudi I. Mozetiču in N. Lauruž za konstruktivne pripombe.

#### Zaključek analize:

Natančna analiza algoritmov za iskanje rezerviranih besed da določeno prednost pomolnim minimalnim ali skoraj minimalnim razpršitvenim funkcijam. Algoritmi z razdelitvijo rezerviranih besed po pruh dveh vrkah je dasovno celo malenkost hitrejši, vendar troši bistveno več spominskega prostora. Pomolne razpršitvene funkcije imajo le to dodatno neprijetnost, da moramo najti oblike razpršitvenih funkcij in da moramo prirediti vrednosti vrkam. To pa je običajno dela nekaj dni.

#### 4. Zaključek

Praktični in teoretični rezultati dajejo prednost uporabi algoritmov z minimalnimi ali skoraj minimalnimi razpršitvenimi funkcijami za iskanje rezerviranih besed. Sama uporaba razpršitvene funkcije je 5 do 10 krat hitrejša od osnovnega binarnega iskanja. Zaradi podatnega branja je v praktičnih primerih prihranek precej manjši, vendar se vedno uspeševanja vreden. Zato je smiselno počitriti obstoječo programsko opremo s pomolno razpršitveno funkcijo.

#### Literatura:

1. R.J.Cichelli: Minimal Perfect Hash Functions made simple, CACM, Vol.23, Num.1, Januar 1980, str. 17-19.
2. G.Jaeckle, G.Osterburg: On Cichelli's Minimal Perfect Hash Function Method, CACM, Vol.23, Num.12, December 1980, str. 729-729.
3. G. Jaeschke: Reciprocal Hashing: A Method For Generating Minimal Perfect Hashing Functions, CACM, Vol.24, Num.12, December 1981, str. 829-833.