

Heširanje

Idealno pretraživanje je kada na osnovu zadatog ključa možemo da direktno pristupimo podatku bez upoređivanja sa drugim ključevima. Ovo bi bilo moguće ukoliko bi obezbedili tabelu koja ima onoliko mesta koliko ima različitih mogućih vrednosti ključeva. Glavni problem je ovde kada je broj mogućih ključeva daleko veći od broja potrebnih ključeva pa je onda iskoriscenje tabele vrlo malo.

Za resenje ovog problema potrebno je pronaci funkciju koja transformise kljuc u ceo broj koji pripada opsegu indeksa u tabeli. Upravo je **hesiranje** tehnika kojom se vrši preslikavanje skupa kljuceva na tabelu znacajno manjih dimenzija. Idealno bi bilo da funkcija za svaki kljuc daje jedinstvenu poziciju. Ta funkcija naziva se **hes funkcija**, a tabela koja se koristi u tom postupku zove se **hes tabela**. Ponekad se ovaj metod hesiranja naziva i tehnikom **rasutog adresiranja ili tehnikom direktnog adresiranja**. Ulaz hes tabele moze pored kljuceva da sadrzi i same zapise ili pokazivace na zapise. Kapacitet ulaza tabele je jedan kljuc ili jedan zapis koji kljuc identifikuje. Ako je velicina hes tabele **n** , onda su indeksi od **0** do **$n-1$** .

Neka je $T[i]$, $0 \leq i < n-1$, hes tabela sa n ulaza i neka ključevi pripadaju nekom skupu mogućih vrednosti S . Broj ključeva je tipično mnogo veći od broja ulaza u tabeli. Neka je h hes funkcija koja vrsi mapiranje ključeva u adresni prostor tabele. Tada funkcija primenjena na neki ključ K daje **matricnu adresu i ključa K , $i = h(K)$** .

42.a) (03.septembar 2008.)

Objasniti i analitički definisati pojam kolizije.

odgovor:

Moguće je da jedna funkcija za dva različita ključa da istu vrednost:

$$h(K_i) = h(K_j), \quad K_i \neq K_j$$

Posto dva ključa ne mogu biti smestena na isto mesto u tabeli, ovaj slučaj predstavlja problem koji se naziva **kolizija**. Ključevi K_i i K_j se nazivaju **sinonimi**, a skup sinonima je **klasa ekvivalencije**. Kolizije je praktično nemoguće izbeći ukoliko ključevi nisu unapred poznati.

Pozeljne osobine za efikasnost pristupa su:

- **uniformnost** – za hes funkciju se kaze da je uniformna na skupu kljuceva S kada za slucajni kljuc $K \in S$ postoji jednaka verovatnoca da se on mapira u bilo koju maticnu adresu:

$$P(i = h(K)) = 1/n, \quad 0 \leq i \leq n-1$$

- **odrzavanje poretka** – sekvencijalan pristup u poretku kljuceva bi se mogao ostvariti hes funkcijama koje odrzavaju poredak i zadovoljavaju uslov:

$$h(K_i) > h(K_j) \text{ za } K_i > K_j$$

Pri primeni tehnike hesiranja treba napraviti dve skoro nezavisne odluke:

- izabrati efikasnu hes funkciju
- izabrati metod za razresavanje problema kolizije

Heš funkcija

Kriterijum za izbor hes funkcije:

- **jednostavnost** – radi brzog izracunavanja
- **uniformnost** – u cilju smanjenja verovatnoce kolizije

U praksi kljucevi mogu biti numericki, alfanumericki, alfabetski. Vecina hes funkcija trazi da kljuc bude iskljucivo numericki zbog aritmetickih operacija nad njim, pa alfabetski i alfanumericki kljucevi treba da se predstave pre primene hes funkcije u numerickom obliku.

Metod ekstrakcije

Ovo je izuzetno jednostavan metod koji se svodi na izvlacenje (ekstrakciju) odredjenog potrebnog broja cifara sa unapred definisanih pozicija.

Ako imamo kljuc koji se na primer sastoji od 7 cifara i imamo tabelu od na primer 1000 (0 do 999) ulaza. Nama su potrebne tri cifre da adresiramo ovu tabelu. Mi sada uzmemo cifre sa nekih definisanih pozicija. Neka su to pozicije 2, 4 i 5. Potom slepimo, formiramo indeks i udjemo u tabelu i to nam je matricna adresa. Ovom metodu je najbolja osobina jednostavnost, a sve ostalo je diskutabilno jer ce sinonimi biti svi oni kljucevi koji imaju iste cifre na ove tri pozicije. Pozeljno je da hes funkcija zavisi od svih znakova kljuca, ali i od njihovih pozicija jer bi se u suprotnom kljucevi koji se mogu dobiti prostom permutacijom znakova mapirali na iste adrese.

Kako biramo hes funkcije?!

Hes funkcije mogu biti:

- **nezavisne od raspodele kljuceva** – ovo je cesca varijanta gde nam nisu poznate verovatnoce pojavljivanja kljuceva
- **zavisne od raspodele kljuceva** – ovo je povoljniji slucaj jer mi tacno znamo koji ce nam kljucevi dolaziti. Onda mozemo da odaberemo uniformniju hes funkciju.

Heš funkcije nezavisne od raspodele ključeva

Metod deljenja

Rezultat hes funkcije je ostatak pri deljenju celobrojnog kljuca nekim brojem n koji je manji ili jednak velicini hes tabele:

$$h(K) = K \bmod n$$

Performanse ovog metoda u smislu smanjenja verovatnce kolizije veoma zavise od pravilnog izbora delioca n .

45. (februar 2012.)

Objasniti kako ne treba odnosno kako treba birati delilac kod metoda deljenja pri hesiranju.

odgovor:

Ne preporučuje se da n bude:

- **parno** – jer se tada parni ključevi mapiraju u parne ulaze, a neparni u neparne ulaze. Ako su ključevi pretežno parni ili neparni, samo polovina ulaza će biti opterećena.
- **stepen broja 2 ili 10** – delilac koji je stepen broja 2, iako omogućava lako izračunavanje ostatka izdvajanjem donjih i bitova ključa, upravo zbog toga nije preporučljiv jer ne zavisi od svih bitova ključa. Delilac kao stepen broja 10 nije pogodan ako su ključevi decimalni brojevi.
- **neprosto sa modulom kongruencije** – u slučaju da su ključevi kongruentni po modulu d treba da se izbegava vrednost n koja nije uzajamno prosta sa d

Preporučuje se da n bude:

- **prost broj ne previše blizu stepena broja 2**

Metod množenja

Metod je zasnovan na množenju ključa pogodno odabranom realnom konstantom sa vrednošću između 0 i 1. Zatim se od proizvoda uzme samo deo iza decimalne tačke, što bi trebalo da predstavlja slučajni broj koji zavisi od svih cifara ključa. Na kraju se dobijeni broj pomnoži sa veličinom tabele i dobije pozicija ulaza u tabelu.

$$h(K) = n(cK \bmod 1), \quad 0 < c < 1$$

Ovde izbor delioca n nije kritičan. Ipak, za performanse je bitan izbor vrednosti konstante c . Neke analize su pokazale da je povoljno birati konstantu kao $c = 0.61803$.

Metod sredine kvadrata

Numericka reprezentacija ključa se pomnoži sama sa sobom, pa se sa fiksnih pozicija iz sredine kvadrata uzme onoliko cifara koliko je potrebno za adresiranje tabele.

$$\begin{array}{r} \text{5894} * \text{5894} \\ \hline \text{23576} \\ \text{53046} \\ \text{47152} \\ \text{29470} \\ \hline \text{34739236} \\ \downarrow \\ \text{39} \end{array}$$

Metod sklapanja

Metod se zasniva na deobi ključa na delove iste dužine. Zatim se ovi delovi sabere, a prenos se ignosrise.

U ovoj metodi hes funkcija zavisi od svih cifara, ali ne i od njihovih pozicija, pa je u mnogim slucajevima efikasnost pod znakom pitanja. Da bi se malo poboljsala uniformnost, koristi se druga varijanta ove metode.

U drugoj koloni je demonstrirana varijanta ovog metoda sa obrtanjem kod koje se redosled cifara u svakom drugom delu invertuje pre primene operacije sabiranja.

$$\begin{array}{r} \text{19} \\ \text{65} \\ \text{30} \\ \text{14} \\ \hline \text{128} \\ \downarrow \\ \text{28} \end{array} \quad \begin{array}{r} \text{19} \\ \text{56} \\ \text{30} \\ \text{41} \\ \hline \text{146} \\ \downarrow \\ \text{46} \end{array}$$

Metod konverzije osnove

Ako je ključ dat u brojnom sistemu sa osnovom p (obicno 2 ili 10), ovaj metod tretira ključ kao broj u osnovi q ($q > p$).

Ako ključ ima vrednost 6154 u dekadnom brojnom sistemu ($p = 10$), kada se prebaci u sistem $q=13$, dobije se:

$$6 * 13^3 + 1 * 13^2 + 5 * 13^1 + 4 * 13^0 = 13420$$

pa se izdvoji potreban broj cifara. Da bi se smanjila verovatnoca kolizije, q se odabira tako da bude uzajamno prosto sa p .

Savrsena hes funkcija

Hes funkcija koja obezbedjuje da nema kolizije u nekom skupu ključeva se naziva savrsena hes funkcija. Ona ima osobinu:

$$h(K_i) \neq h(K_j), \text{ za svako } i \neq j \text{ gde } K_i, K_j \in S$$

Savrsenu hes funkciju nije lako pronaci, pogotovo za dinamički skup. Sto je hes tabela veca u odnosu na broj ključeva, savrsenu hes funkciju je lakse naci. Idealno bi bilo imati savrsenu hes funkciju za n ključeva u tabeli sa n ulaza. Takva savrsena hes funkcija se naziva **minimalnom**.

Heš funkcije zavisne od raspodele ključeva

Kada je skup ključeva unapred poznat, onda je moguće naći efikasnu hash funkciju koja će smanjiti broj kolizija.

////////////////////

40. (februar 2007.)

Objasniti metodo hesiranja koja koristi **analizu cifara**.

odgovor:

Metod analize cijfara

Neka su kljucevi poznati unapred. Mi sada zelimo da primenimo prost metod ekstrakcije, tj. da hes formiramo od cifara sa unapred izabranih pozicija, ali treba pravilno izabrati pozicije. Prvo se analizira skup numerickih vrednosti svih kljuceva, pa se napravi tabela koja za svaku poziciju kljucja daje broj prijavljivanja pojedinih cifara na toj poziciji u svim kljucevima. Zatim se selektuje onoliko pozicija koliko je cifara potrebno da se adresira hes tabela tako sto se odaberu kolone tabele koje pokazuju najmanje varijacije pojavljivanja razlicitih cifara.

Razresavanje kolizije

Kolizija predstavlja slucaj kada hes funkcija mapira dva ili vise razlicitih kljuceva na istu maticnu adresu. Broj kolizija se najjednostavnije moze smanjiti povecanjem velicine tabele. Tako se poboljsavaju performanse hesiranja, ali po cenu dodatnog prostora i njegovog slabijeg iskoriscenja.

Pri kvantitativnom određivanju performansi najcesce se koristi parameter koji se naziva **popunjenost tabele α (load factor)**. Ovaj parameter predstavlja odnos broja zauzetih ulaza u tabeli i velicine tabele.

Za razresavanje kolizije tradicionalno se koriste dva načina:

- **otvoreno adresiranje** – kolizije se razresavaju u okviru tabele
- **ulancavanje** – kolizije se razresavaju van tabele

Otvoreno adresiranje

*Otvoreno adresiranje se svodi na pronalazjenje neke druge lokacije u hes tabeli različite od maticne adrese kada se javi kolizija. Osnovna ideja je formulisati pravilo koje za svaki ključ određuje **ispitni niz** kao niz adresa u tabeli koje se proveravaju pri umetanju novog ključa. Ako treba da se umetne novi ključ, a maticna adresa je zauzeta, onda se pokušava sa sledećom adresom u ispitnom nizu, sve dok se ne nađe slobodna lokacija.*

Generisanje ispitnog niza za dati ključ naziva se **ponovnim hesiranjem** jer se na osnovu zauzete adrese izracunava neka druga adresa sve dok se ne dodje do slaganja ključa pri pretraživanju ili do prve slobodne lokacije pri umetanju.

Pozeljno je da ispitni niz bude neka permutacija od skupa adresa $(0 \dots n-1)$ u tabeli sto omogućava da se u najgorem slucaju kolizije razrese posle n koraka.

Varijante otvorenog adresiranja su:

- linearno pretrazivanje
- slučajno pretrazivanje
- kvadratno pretrazivanje
- dvostruko hesiranje

////////////////////

44. (februar 2007.)

Dati pseudokod i objasniti opste algoritme umetanja i pretrazivanja ključa u hes tabeli kod **tehnika otvorenog adresiranja**. Kakav se problem javlja kod brisanja i zašto?

odgovor:

```

HASH-INSERT( $T, K$ )
 $i = 0$ 
repeat
     $j = h(K)$ 
    if ( $T[j] = \text{empty}$ ) then
         $T[j] = K$ 
        return  $j$ 
    else
         $i = i + 1$ 
end_if
until  $i = n$ 
ERROR(Tabela cheia)

```

```

HASH-SEARCH( $T, K$ )
 $i = 0$ 
repeat
     $j = h(K)$ 
    if ( $T[j] = K$ ) then
        return  $j$ 
    else
         $i = i + 1$ 
end_if
until ( $T$  is empty) or ( $i = n$ )
return empty

```


Umetanje ključa K u tabelu T se realizuje funkcijom **$HASH-INSERT(T, K)$** . Inicijalno je broj pokušaja pri umetanju jednak nula ($i=0$). Imamo jednu petlju koja generise ispitni niz.

U prvom prolazu pravimo matičnu adresu ($j = h(K)$). Potom proveravamo da li je lokacija prazna ($T[j] = \text{empty}$).

Ako je lokacija prazna, smestamo ključ u istu ($T[j] = K$). Funkcija vraća adresu u tabeli na kojoj je ključ smesten (**$return j$**). Ako lokacija nije prazna, idemo u sledeći pokušaj ($i = i + 1$).

Ukoliko se iscrpi citav niz (**$until i = n$**), a ne nađe se slobodna lokacija, javlja se poruka da je tabela puna (**$ERROR (Tabela puna)$**).

Pretraživanje ključa K u tabeli T se realizuje funkcijom **$HASH-SEARCH(T, K)$** . Postupak je sličan prethodnom. Hesiramo ključ i dodjemo do adrese do koje bi dosli i da umecemo. Proverimo da li je ključ tamo ($T[j] = K$), pa ako jeste, funkcija vraća adresu na kojoj je ključ pronađen (**$return j$**).

Ako nismo pronašli ključ, idemo u rihesiranje ($i = i + 1$). Bitno je istaci da se generisanje ispitnog niza obavlja na isti način prilikom umetanja i prilikom pretraživanja.

Neuspesno pretraživanje se proglašava ako u ispitnom nizu dodjemo do prazne lokacije ili ako obidjemo sve ulaze tabele (**$T[j] = \text{empty or } (i = n)$**), a ključ nije pronađen. Tada funkcija vraća **$return \text{empty}$** .

Problem brisanja se svodi na činjenicu da se prilikom brisanja ključ ne može jednostavno ukloniti iz tabele i taj ulaz proglašiti slobodnim, jer se može prekinuti ispitni niz do nekih drugih ključeva.

Ovaj problem se rešava tako što se na mestu obrisano ključa ostavi poseban kod (na primer, *deleted*) koji je različit od vrednosti svih ključeva, ali i od vrednosti *empty*. Ovo je takozvani **metod "poluslobodne" lokacije**. Sada se pri pretraživanju ispitni niz ne prekida kada se dodje do lokacije *deleted*, već se preskace kao da je zauzeta i ispitni niz dalje generise. Pri umetanju se ovakva lokacija tretira kao slobodna, tako da ona može da se iskoristi za smestanje novog ključa.

Mana ovog postupka je ta što jednom zauzete lokacije više nikada ne postaju slobodne, što rezultira lošim performansama pretraživanja. Zato se ovaj metod koristi kada su brisanja u tabeli redja.

Postoji i **metod selektivnog pomeranja zapisa** kod kojeg se oslobodjena pozicija popuni. To je u principu složena operacija.

////////////////////////////////////

Linearno pretrazivanje

42.d) (03.septembar 2008.)

Objasniti i analiticki definisati pojam linearnog pretrazivanja.

odgovor:

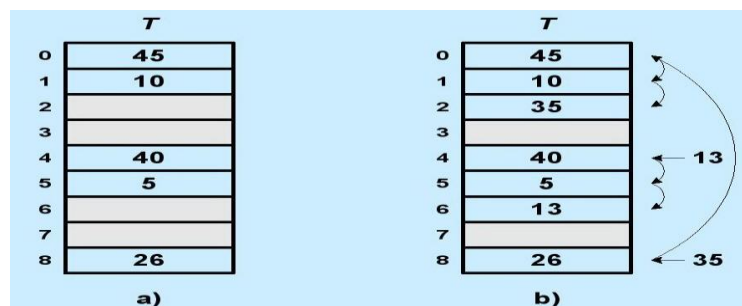
Ovaj metod pri koliziji generise ispitni niz sukcesivnih lokacija, a kad dodje do kraja tabele, vraca se na pocetak. Postupak ponovnog hesiranja se moze predstaviti na jedan od dva nacina:

$$h_i(K) = (h_0(K) + i) \bmod n, \quad i = 1, 2, \dots, n-1$$

$$h_i(K) = (h_{i-1}(K) + 1) \bmod n$$

Linearno pretrazivanje u slucaju kolizije pocinje sa sekvencijalnim trazenjem od maticne adrese i tezi da postavi kljuc sto blize svojoj maticnoj adresi.

Primer:



Neka je data prazna tabela sa 9 ulaza i hes funkcija $h(K) = K \bmod 9$. Kljucevi 5, 10, 26, 40, 45 se tada umecu bez kolizije na svoje maticne adrese (*slika a*)).

Medjutim kada dodje kljuc 13, njegova maticna adresa 4 je zauzeta, pa se pokusava na prvu sledecu, ali je i ona zauzeta jer se tu nalazi kljuc 5. Proba se sa sledecom adresom, ona je slobodna , tako da kljuc 13 popunjava ulaz 6.

Prilikom umetanja kljuca 35, njegova maticna adresa 8 je zauzeta, pa se pokusava na adrese 0, 1 i tek je lokacija 2 slobodna.

Selektivno pomeranje pri brisanju kod linearnog pretraživanja:

Kod linearnog pretraživanja je moguće premještanje ključeva posle brisanja sa ciljem da se popuni obrisano mjesto, a da se ne prekine ispitni niz ka ključevima na narednim adresama.

Ako je brisanje izvršeno na lokaciji i , tada se traži prva sledeća prazna lokacija j po pravilu linearnog pretraživanja. Zatim se počevši od adrese $i+1$, izračunavaju matične adrese za ključeve koji se nalaze između lokacija i i j .

Neka se na nekoj takvoj lokaciji l nalazi ključ K . Ako je $\mathbf{r} = \mathbf{h}(K)$ matična adresa tog ključa i ako se ona ne nalazi između i i j ($\mathbf{r} < i$ ili $\mathbf{r} > j$), onda je moguće izvršiti premještanje ključa sa tekuće lokacije na lokaciju i i time osloboditi tekucu poziciju.

Ako se r nalazi izmedju i i j , onda se kljuc ne moze premestiti jer bi se time on postavio ispred svoje maticne adrese sa koje pocinje pretrazivanje na njega.

Posto je sada prazna lokacija presla na adresu l , postupak se zatim sukcesivno ponavlja sa lokacijama izmedju l i j , sa ciljem da ova slobodna pozicija “izadje” na adresu $j-l$.

Ovakvim premestanjem se popravljaju performanse pretrazivanja.

Primarno grupisanje

////////////////////

42.b) (03.septembar 2008.)

Objasniti i analitički definisati pojam **primarnog grupisanja** (i kako se ono izbegava).

odgovor:

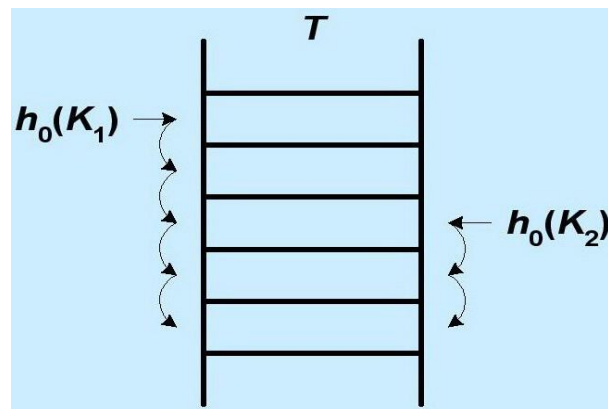
Glavna prednost linearnog pretraživanja je jednostavnost, ali je glavni nedostatak pojava **primarno grupisanje** koja se obično javlja kod veće popunjenosti tabele.

Analiticki se definiše kao:

$$h_{i+r}(K_1) = h_r(K_2) \quad \text{za } r = 0, 1, 2, \dots$$

Ako imamo ispitni niz od nekog ključa \mathbf{K}_I pa u nultom pokusaju dodjemo do matične adrese, u prvom pokusaju dodjemo do neke sledeće i tako redom. Ako se sada desi da matična

adresa nekog drugog ključa K_2 upadne na neko mesto u ispitni niz prvog ključa i od tog trenutka se oni poklapaju, onda se ta pojava naziva **primarno grupisanje**.



Ova pojava nije dobra jer će sinonimi od K_1 i K_2 završavati na istim mestima i onda će se ispitni niz produžavati pa će se stvarati grupe zauzetih lokacija i neće se ključevi razbacivati po samoj tabeli.

Kako rešiti problem primarnog grupisanja?!

Jednostavan način da se ovo izbegne je da hes funkcija bude zavisna i od broja pokusaja:

$$h_i(K) = (h_{i-1}(K) + i) \bmod n, \quad i = 1, 2, \dots$$

////////////////////////////////////

42.e) (03.septembar 2008.)

Objasniti i analitički definisati pojam **kvadratnog pretraživanja** i objasniti njegove prednosti.

odgovor:

Jos bolji način da se izbegne primarno grupisanje je primena tehnike **kvadratnog pretraživanja** kod koje se ispitni niz generise preko kvadrata broja neuspesnih pokusaja:

$$h_i(K) = (h_0(K) + i^2) \bmod n, \quad i = 1, 2, \dots, n-1$$

Ovaj niz se može generisati i bez dodatnih operacija kvadriranja ako se ima u vidu da je:

Dvostruko heširanje

42.f) (03.septembar 2008.)

Objasniti i analitički definisati pojam *dvostrukog heširanja*.

41. (14.februar 2007.)

Objasniti da li se primarno i sekundarno grupisanje izbegava kod dvostrukog heširanja.

odgovori:

Ovaj metod se zasniva na primeni dve nezavisne hes funkcije: primarne $h(K)$ i sekundarne $g(K)$:

$$h_i(K) = (h_{i-1}(K) + g(K)) \bmod n$$

ili

$$h_i(K) = (h_0(K) + i g(K)) \bmod n, \quad i = 1, 2, \dots, n-1$$

Na zadati ključ se prvo primenjuje primarna hes funkcija, pa ako je došlo do kolizije, na vrednost ključa se primenjuje i sekundarna hes funkcija koja onda izračunava pomeraj u ispitnom nizu.

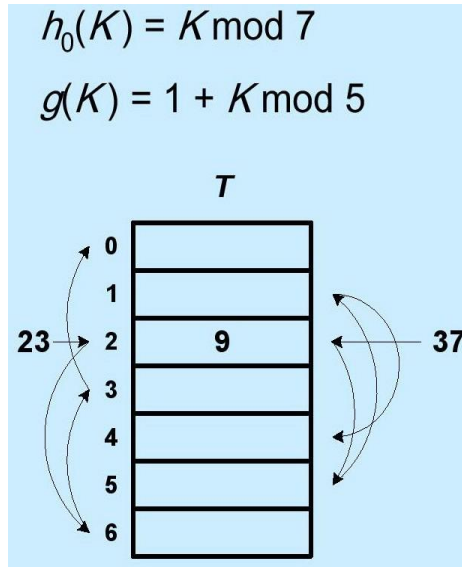
Dvostruko heširanje veoma smanjuje verovatnoću sekundarnog grupisanja za nezavisne funkcije h i g (ne otklanja skroz).

Sekundarna hes funkcija mora da bude pažljivo izabrana da bi ispitni niz sadržao sve adrese. Zato $g(K)$ mora da daje rezultat u opsegu 1 do $n-1$ i da bude uzajamno prost sa n .

Knuth predlaže

- ✓ $h_0(K) = K \bmod n$
- ✓ $g(K) = 1 + K \bmod (n-2)$
- ✓ n i $n-2$ prosti brojevi

Primer dvostrukog heširanja:



Ključevi 23 i 37 su sinonimi sa istom matricnom adresom 2 jer je:

- $23 \bmod 7 = 2$
- $37 \bmod 7 = 2$

Kada je ona zauzeta (u ovom primeru ključem 9), dvostruko hesiranje generise dva razlicita ispitna niza. To se radi po principu:

(“dobijen ulaz koji nije slobodan” + $g(K)$) mod n (gledas n iz primarne heš funkcije, u nasem primeru je to 7) :

Ispitni niz za ključ 23 (pored lokacije 2) cine:

- $(2 + 1 + 23 \bmod 5) \bmod 7 = 6 \bmod 7 = 6$
- $(6 + 1 + 23 \bmod 5) \bmod 7 = 10 \bmod 7 = 3$
- $(3 + 1 + 23 \bmod 5) \bmod 7 = 7 \bmod 7 = 0$ itd.

Ispitni niz za ključ 37 (pored lokacije 2) cine:

- $(2 + 1 + 37 \bmod 5) \bmod 7 = 5 \bmod 7 = 5$
- $(5 + 1 + 37 \bmod 5) \bmod 7 = 8 \bmod 7 = 1$
- $(1 + 1 + 37 \bmod 5) \bmod 7 = 4 \bmod 7 = 4$ itd.

Poboljsanja metoda otvorenog adresiranja

Kada je zauzetost tabele velika broj poredjenja pri pretrazivanju raste, sto znatno usporava ovu operaciju. Zato uvodimo dve metode poboljsanja gde **metoda uredjene tabele** smanjuje poredjenja pri neuspesnom pretrazivanju, a **Brent-ov metod** smanjuje broj poredjenja pri uspesnom pretrazivanju.

Metod uredjene tabele

Ovaj metod pri umetanju odrzava skup kljuceva koji se mapiraju na istu maticnu adresu uredjenim po opadajucim vrednostima kljuceva. Zbog toga, kada se u operaciji umetanja novog kljuca K pri trazenju praznog mesta u ispitnom nizu dodje do prve lokacije na kojoj se nalazi kljuc $K' < K$, kljuc K se umece na njegovo mesto, a zatim se nastavlja sa trazenjem mesta za K na isti nacin.

Kada se koristi?!

Moze se primeniti na sve tehnike kod kojih se generisanje sledece adrese u ispitnom nizu vrsi samo na osnovu prethodne adrese i vrednosti kljuca, ali ne i na osnovu prethodnog broja pokusaja.

Performanse

Znacaj rezultata ovog metoda je sto izjednacava broj poredjenja pri uspesnom i neuspesnom pretrazivanju, ali kviri performanse umetanja, pa je bolji za staticke tabele koje cesto pretrazuju.

Primer

$$h_i(K) = (K + 2i) \bmod n$$

a)

	T
0	
1	55
2	
3	28
4	
5	19
6	
7	
8	

b)

	T
0	
1	55
2	
3	37
4	
5	28
6	
7	19
8	

Na slici a) je prikazano stanje nakon umetanja kljuceva 55, 19 i 28 tim redom. Svi imaju istu maticnu adresu 1. Kada se umece kljuc 37 (*slika b*)), koji je sinonim ovih kljuceva, on se ne stavlja na kraj ispitnog niza (1, 3, 5, 7...) vec se umece ispred kljuca 28. Na taj nacin se odrzava skup kljuceva koji se mapiraju na istu maticnu adresu po opadajucim vrednostima kljuceva. Ako bi nakon ovoga, pretrazivali tabelu na kljuc 46, ne moramo proveriti sve lokacije vec odmah posle provere na adresi 3 ($37 < 46$) utvrdjuje se da ovaj kljuc nije u tabeli. Dakle, dobitak je taj sto ne moramo da idemo do kraja ispitnog niza.

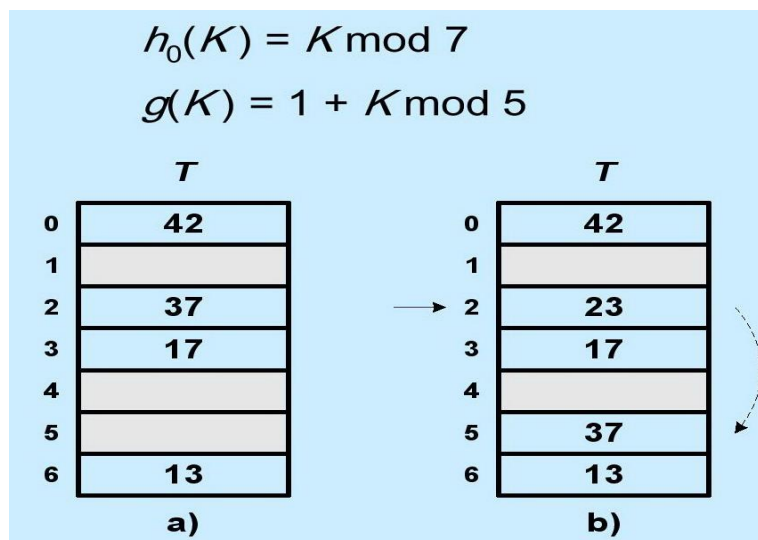
Brent-ov metod

Ovaj metod smanjuje broj poredjenja pri uspesnom pretrazivanju, a nema uticaja na neuspesno pretrazivanje. Zasniva se na preuredjenju hes tabele i koristi se pri dvostrukom hesiranju. *Osnovna ideja je da se premestanjem nekog vec umetnutog kljuc oslobodi mesto za novi kljuc, a da ukupan prosecan broj poredjenja pri trazenju jednog i drugog kljuc bude smanjen.*

Kada se koristi?!

Ovaj metod je pogodan za staticke tabele gde nema mnogo umetanja.

Primer



Pocetno stanje tabele dato je na slici a). Mi sada zelimo da umetnemo kljuc 23. Njegov ispitni niz cine lokacije 2, 6, 3, 0 i sve su one popunjene. Kljuc 23 bi mogao da se postavi tek na prvu slobodnu lokaciju, a to je adresa 4. U ovoj varijanti bi za uspesno pretrazivanje kljuc 23 bilo potrebno 5 poredjenja (2, 6, 3, 0, 4), a za kljuc 37 samo jedno predjenje, a u proseku 3 poredjenja po kljucu.

Ako bi primenili **Brent-ov** metod pa kljuc 23 postavili na maticnu adresu 2 ($23 \bmod 7 = 2$), a kljuc 37 pomerili na prvu sledecu lokaciju koja je slobodna iz njegovog ispitnog niza (2, 5, 1...), a to je adresa 5, tada bi pri uspesnom pretrazivanju ova dva kljuc bilo potrebno ukupno 3 poredjenja (jedno za kljuc 23 i dva za kljuc 37), a u proseku 1.5 poredjenja.

Ulančavanje

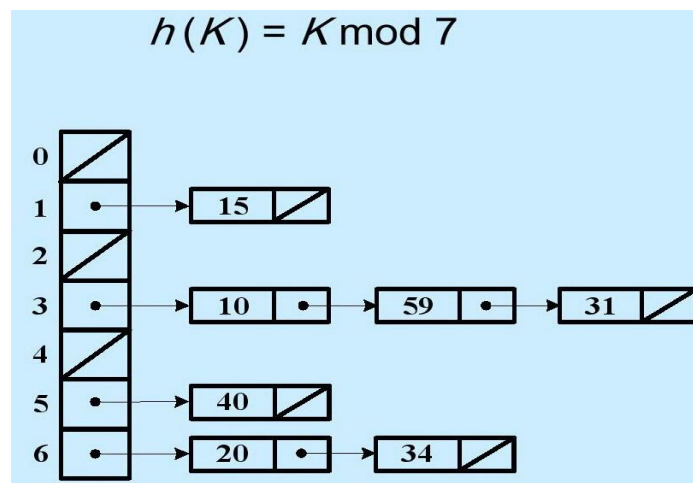
Tehnike otvorenog adresiranja imaju dva nedostatka. Prvi problem nastaje u primenama gde ima dosta brisanja, jer se zauzete lokacije ne oslobadjaju, a premestanje može biti vrlo složeno, ako je uopšte moguće.

Drugi problem je neefikasno pretraživanje kod tabele sa većom popunjenosti, jer se u ispitnom nizu mogu naći i adrese zauzete ključevima koji nisu sinonimi traženog ključa. Pomenuti problemi se efikasno rešavaju drugom tehnikom koja se zove **ulančavanje**.

Odvojeno ulančavanje

Sustina ove tehnike je da se sinonimi ulančaju u liste čija se zaglavlja nalaze na matičnim adresama. Zauzeti ulazni hes tabele sadrži samo pokazivač na listu sinonima koji odgovaraju toj matičnoj adresi. Pošto svaka lista odgovara samo jednoj klasi ekvivalencije, ova varijanta ulančavanja se naziva **odvojeno ulančavanje**.

Primer



Za ključ 15 matična adresa je 1, a za ključeve 10, 59 i 31 matična adresa je 3 pa su oni ulančani. Za ključ 40 matična adresa je 5, a za ključeve 20 i 34 matična adresa je 6 pa su oni ulančani.

////////////////////////////////////

43. (24.septembar 2008.)

Koje su prednosti, a koji nedostaci odvojenog ulancavanja u odnosu na otvoreno adresiranje?

odgovor:

Prednosti:

Operacije sa ključevima kod tehnike ulancava su jednostavnije nego kod otvorenog adresiranja.

Pretraživanje na zadati ključ se svodi na prolazak kroz ulancanu listu čije se zaglavlje nalazi na matičnoj adresi. Ako se citava lista prodje, a ključ se ne nađe, pretraživanje je neuspješno.

Umetanje novog ključa može da se vrši na početak liste, ako se zna da ključ nije u tabeli, ili na kraj liste, ako treba prvo proveriti da li je ključ unesen, jer se tu završava neuspješno pretraživanje. Ako se lista sinonima drži uređena po vrednosti ključa, tada je pretraživanje efikasnije, ali je umetanje manje efikasno jer se ključ mora postaviti na mesto koje mu po poretku odgovara.

Brisanje je veoma jednostavno jer se svodi na prevezivanje liste. Zbog brisanja je pogodno da lista bude dvostruko ulancana.

Može se očekivati da prosečne dužine lista budu mala, pa su operacije sa hes tabelom brze. Kada bi liste bile dugacke, onda bi bilo efikasnije da se umesto liste koristi binarno stablo pretraživanja.

Dodatna prednost je i to što **tehnika ne ograničava broj ključeva** kao kod otvorenog adresiranja gde je broj ključeva ograničen veličinom kontinualne tabele.

Nedostaci:

Sve ove prednosti se plaćaju dodatnim prostorom za pokazivace i većim brojem praznih ulaza u hes tabeli.

////////////////////////////////////

Objedinjeno ulancavanje

Sto je veca hes tabela, liste su u proseku krace, a performanse bolje. Sa povecanjem hes tabele, raste i broj nezauzetih ulaza u hes tabeli. Zato se javlja ideja da se taj prostor iskoristi za smestaj elemenata liste, sto poboljsava efikasnost koriscenja prostora sa dva aspekta: koristi nezauzete ulaze hes memorije i izbegava dodatni prostor za smestaj clanova liste. Ovakva varijanta ulancavanja se naziva **standardno objedinjeno hesiranje**.

Ovaj metod predstavlja hibridno resenje izmedju otvorenog adresiranja i odvojenog ulancavanja, jer kljuceve smesta samo u osnovnoj tabeli, kao otvoreno adresiranje, ali ne koristi ponovno hesiranje u slucaju kolizije vec direktne pokazivace na sledecu adresu u ispitnom nizu.

Svaki ulaz hes tabele treba da pored kljuka key sadrzi i polje next koje sluzi za cuvanje adrese ulaza u tabeli gde se nalazi eventualni naredni clan liste sinonima. U pocetnom stanju u svim ulazima polje kljuka ima pocetno stanje empty, a polje next je -1, sto je analogno praznom pokazivacu.

Objedinjeno ulancavanje bolje koristi prostor u odnosu na odvojeno ulancavanje, ali su performanse losije. Nepovoljan efekat ove metode je sto dopusta objedinjavanje lista za vise klase ekvivalencije. U primeru koji sledi kljucevi 42, 62, 78 i 88 se nalaze u jednoj listi, a to su kljucevi iz dve klase ekvivalencije za maticne adrese 2 i 8. Spajanje se desilo kada je kljuc 62 postavljen na lokaciju 8 koja predstavlja maticnu adresu za kljuceve 88 i 78, jer je njegova maticna adresa 2 vec bila popunjena kljucem 42.

Primer

0	50	-1
1	-	-1
2	42	8
3	-	-1
free → 4	78	-1
5	25	-1
6	19	-1
7	88	4
8	62	7
9	9	6

Hes tabela ima 10 ulaza, a hes funkcija je oblika $h(K) = K \bmod 10$.

U tabelu se redom umecu kljucevi : 42, 9, 25, 62, 88, 50, 19, 78.

- $h(42) = 42 \bmod 10 = 2$; ulaz 2: 42; next 2: -1;
- $h(9) = 9 \bmod 10 = 9$; ulaz 9: 9; next 9: -1;
- $h(25) = 25 \bmod 10 = 5$; ulaz 5: 25; next 5: -1;
- $h(62) = 62 \bmod 10 = 2$; ulaz 2 je zauzet pa kljuc smestamo na prvu slobodnu lokaciju sa kraja na koju ukazuje **pokazivac free**, a to je ulaz 8; next 2: 8;
- $h(88) = 88 \bmod 10 = 8$; ulaz 8 je zauzet pa kljuc smestamo na prvu slobodnu lokaciju sa kraja na koju ukazuje, a to je ulaz 7; next 8: 7;
- $h(50) = 50 \bmod 10 = 0$; ulaz 0: 50; next 0: -1;
- $h(19) = 19 \bmod 10 = 9$; ulaz 9 je zauzet pa kljuc smestamo na prvu slobodnu lokaciju sa kraja, a to je ulaz 6; next 9: 6;
- $h(78) = 78 \bmod 10 = 8$; ulaz 8 je zauzet pa kljuc smestamo na prvu slobodnu lokaciju sa kraja, a to je ulaz 4; next 7: 4;

Komparativno poredjenje tehnika unutrasnjeg hesiranja

U cilju poredjenja performansi, u narednoj tabeli su prikazani brojevi pokusaja pri uspesnom i neuspesnom pretrazivanju za odredjene vrednosti popunjenosti (α) hes tabele kod:

- linearnog pretrazivanja (LP)
- dvostrukog hesiranja (DH)
- odvojenog ulancavanja (SC)

	S			U		
α	LP	DH	SC	LP	DH	SC
0.1	1.056	1.054	1.050	1.118	1.111	1.005
0.2	1.125	1.116	1.100	1.281	1.250	1.019
0.3	1.214	1.189	1.150	1.520	1.429	1.041
0.4	1.333	1.277	1.200	1.889	1.667	1.070
0.5	1.500	1.386	1.250	2.500	2.000	1.107
0.6	1.750	1.527	1.300	3.625	2.500	1.149
0.7	2.167	1.720	1.350	6.060	3.333	1.197
0.8	3.000	2.012	1.400	13.000	5.000	1.249
0.9	5.500	2.558	1.450	50.500	10.000	1.307
0.95	10.500	3.153	1.475	200.50	20.000	1.337

Za manje popunjenosti tabele, do 50%, performanse ovih metoda su dosta bliske, pogotovo pri uspesnom pretrazivanju.

Razlike u performansama su izuzetno vidljive kada je tabela prilično popunjena (oko 95%). Tehnika odvojenog ulancavanja pokazuje najbolje performanse jer se može očekivati da se u 1.5 poredjenja nadje ključ. Očigledno je da performanse ove metode najmanje variraju sa povećanjem popunjenosti. Mana tehnike odvojenog ulancavanja je potreba za dodatnim prostorom, ima više nezauzetih ulaza u memoriji, a i pokazivaci predstavljaju dodatnu cenu.

Kako popunjenost raste, linearno pretrazivanje pokazuje najslabije performanse. Ipak, valja podsetiti da otvoreno adresiranje za razresavanje kolizije koristi samo osnovnu tabelu, dok odvojeno ulancavanje koristi dinamički prostor van tabele.

Sto se tice izbora hes funkcije, ako raspodela kljuceva nije unapred poznata, metod deljenja se pokazao najprakticnijim.

Kod otvorenog adresiranja imamo dodatni problem što velicina hes tabele ogranicava broj kljuceva. Ako se tabela napuni, a nailaze i dalje kljucevi, resenje moze da bude u povecanju tabele za odredjeni procenat i ponovnom hesiranju svih kljuceva.

Spoljašnje heširanje

Objasniti glavne specifičnosti postupka spoljasnjeg hesiranja u odnosu na unutrašnje.

Hesiranje se kao tehnika pretraživanja koristi i za pristup podacima na spoljanim memorijama, najcesce u datotekama sa direktnim pristupom.

Pristup disku je daleko najskuplja operacija pa se pri spoljasmjem hesiranju mnogo vise vodi racuna o vremenu nego o prostoru pogotovo kad se zna da je prostor na disku veoma jeftin. Zato je jako vazno smanjiti broj kolizija i pristupa disku, cak i na racun nesto slabijeg iskoriscenja prostora.

Ukoliko je baket pun, onda treba da se primeni neka tehnika za razresavanje kolizije jer ključ ne može da stane u svoj matični baket pa mu treba naći mesto u nekom drugom baketu.

Traženje zapisa sa datim ključem počinje učitavanjem matičnog baketa koji se zatim pretražuje u memoriji. Ako je **b** manje, obično se isplati sekvencijalno pretraživanje. Za veće **b**, može da se primeni i binarno pretraživanje ako su ključevi u baketu uređeni. Ako se ključ ne pronađe u matičnom baketu, učitava se drugi baket po istom principu koji je primenjen za razresavanje kolizije kod umetanja. Neuspesno pretraživanje se proglašava ako baket nije pun, a ključ se ne nađe u njemu.

Kod **brisanja**, ukoliko je baket pun, ne možemo baš slobodno da brisemo ključ iz njega. Moramo da ga eventualno dopunimo sa ključem iz sledeće lokacije u ispitnom nizu. Ako baket nije pun, možemo slobodno da brisemo.

Kolizije se kod spoljasnog hesiranja mogu razresavati primenom istih tehnika kao kod unutrašnjeg, a to su otvoreno adresiranje ili ulancavanje.

Analize pokazuju da za manji kapacitet baketa imamo sličnu situaciju kao kod unutrašnjeg hesiranja. Tada je linearno pretraživanje najslabije. Dvostruko hesiranje ima najbolje performanse za iole veće bakete. Mada su generalno gledano performanse prilično ujednačene za veće kapacitete baketa.

Treba voditi računa o činjenici da pored broja citanja baketa, važna je i njihova relativna pozicija zbog vremena pozicioniranja glava na disku. Svakako je brže kada se citaju dva fizički susedna baketa nego ako su oni fizički udaljeni. Ova činjenica istice u prvi plan **linearno pretraživanje**, koje kod kolizije pokušava sa sledećim, susednim baketom.

Kod **odvojenog ulancavanja**, kolizije se razresavaju u posebnoj oblasti prepunjenja gde se smestaju zapisi koji ne stanu u matični baket. Ova oblast može da se organizuje u posebne bakete za svaku listu, ali i ne mora, jer nije previše verovatno da će posle prepunjenja nekog matičnog baketa biti dovoljno ključeva da ispune još jedan baket. Zato se u oblasti prepunjenja mogu ulancavati individualni zapisi u okviru zajedničkih baketa.

Objedinjeno ulancavanje ne koristi posebnu oblast popunjenja već kolizije razrešava ulancavanjem u baketima u primarnoj oblasti koji nisu puni. U tom cilju se održava lista baketa koji imaju slobodnih mesta. Kada se baket napuni, onda se on isključuje iz ove liste.

Nemogućnost sekvencijalnog pristupa zapisima po uređenom poretku njihovih ključeva prisutna kod metoda unutrašnjeg hesiranja, karakterise i spoljasnje hesiranje.

////////////////////////////////////

Fleksibilne tehnike spoljašnjeg heširanja

Datoteke su relativno dugovečne strukture, ali tokom svog trajanja one mogu da se povećavaju ili smanjuju na nepredvidljiv način. U takvim uslovima su metode spoljašnjeg heširanja nefleksibilne. Velicina heš tabele se projektuje prema nekom predviđenom broju zapisa, sa namerom da se postigne određena popunjenost, a time i željene performanse. Međutim, ako tokom rada broj zapisa poraste mnogo iznad predviđenog broja, popunjenost raste, a performanse se pogoršavaju. Ako broj zapisa dosta opadne, onda se prostor neefikasno koristi budući da je velicina heš tabele fiksna.

Prema tome, poželjno bi bilo učiniti spoljašnje heširanje fleksibilnijim. U tom cilju su razvijene tehnike ***dinamičkog, proširljivog i linearnog heširanja***. Mi ćemo na ovom kursu govoriti o dinamičkom i proširljivom heširanju.

Dinamičko heširanja

12. nedelja, 3. čas, [14:15 – 22:35]

Proširljivo heširanja

12. nedelja, 3. čas, [22:35 – 33:40]