



BM 307 Dosya Organizasyonu (File Organization)

Hazırlayan: M.Ali Akcayol
Gazi Üniversitesi
Bilgisayar Mühendisliği Bölümü



Konular

- Direct File Organization
 - Computed Chaining
 - Comparison of Collision Resolution Methods
 - Perfect Hashing
 - Cichelli's Algorithm



Direct File Organization

Computed Chaining

- Link kullanarak çakışma çözümü yapan metodlar (coalesced hashing) ve link kullanmadan çözüm yapan metodlar (progressive overflow, linear quotient, Brent's method, binary tree)
- Link kullanan metodlar ek yer gerektirmekte ancak daha iyi performansa sahiptir
- Link kullanmayan metodlar ek yer gerektirmemekte ancak daha kötü performansa sahiptir
- Computed chaining metodunda link adresleri kayıt edilmezler hesaplanarak bulunurlar
- Probe chain içerisindeki elemanların tam adresi yerine ofset adresleri saklanır
- Birkaç bit kullanılarak adres bilgisine ulaşılabilir



Direct File Organization

Computed Chaining

- Coalesced chain oluşmaması için eklenen kaydın home adresine yazılmış kayıtlar yeni bir adrese taşınır
- Böylece bir probe chain içindeki tüm kayıtlar aynı home adresine sahiptir
- Bir kayıt için probe sayısı o kaydın probe chain içindeki pozisyonuna eşittir
- Computed chain bir sonraki adresi hesaplarırken eklenen veya okunacak kaydı değil o adresteki kaydı kullanır

| | Record | Pseudo-link |
|----------|--------|-------------|
| <i>r</i> | A | $i(A)$ |
| | | |
| <i>s</i> | B | $i(B)$ |
| | | |
| <i>t</i> | C | Λ |
| | | |

Direct File Organization

Computed Chaining

Örnek

- 27, 18, 29, 28, 39, 13, 16, ve 38, 53
- $\text{Hash}(\text{key}) = \text{key} \bmod 11$
- $i(\text{key}) = \text{Quotient}(\text{Key} / 11) \bmod 11$

nof (number of offsets,
pseudolink value)

| | Key | nof | | Key | nof | | Key | nof |
|----|-----|-----|----|-----|-----|----|-----|-----|
| 0 | | | 0 | | | 0 | | |
| 1 | | | 1 | | | 1 | | |
| 2 | | | 2 | | | 2 | 13 | |
| 3 | | | 3 | | | 3 | | |
| 4 | | | 4 | | | 4 | | |
| 5 | 27 | | 5 | 27 | | 5 | 27 | 2 |
| 6 | | | 6 | 28 | 2 | 6 | 28 | 2 |
| 7 | 18 | 1 | 7 | 18 | 1 | 7 | 18 | 1 |
| 8 | 29 | | 8 | 29 | | 8 | 29 | |
| 9 | | | 9 | | | 9 | 16 | |
| 10 | | | 10 | 39 | | 10 | 39 | |

Direct File Organization

Computed Chaining

Değerlendirme

- İlk 7 kayıt için ortalama probe sayısı 1.4 tür.
- Bu sayı binary tree için 1.7 ve LISCH için 1.4 olmuştur.
- LISCH ile aynı performansa sahiptir ancak pseudolink için gereken alan LISCH deki link için kullanılan alandan daha azdır.
- Pseudolink için 2 bit gerekmektedir. Pointer'la tam adres gösterimi için 3 veya 4 byte gerekmektedir.
- En büyük adres 10 olsa bile en fazla 4 bit kullanılması yeterlidir.

Direct File Organization

Computed Chaining

Örnek (devam)

- 38 ve 53 eklendiğinde son durum aşağıdaki gibi olur.
- $\text{Hash}(\text{key}) = \text{key} \bmod 11$
- $i(\text{key}) = \text{Quotient}(\text{Key} / 11) \bmod 11$

| | Key | nof | Key | nof | |
|----|-----|-----|-----|-----|---|
| 0 | 38 | | 0 | 16 | 1 |
| 1 | | | 1 | 38 | |
| 2 | 13 | | 2 | 13 | |
| 3 | | | 3 | | |
| 4 | | | 4 | | |
| 5 | 27 | 2 | 5 | 27 | 3 |
| 6 | 28 | 2 | 6 | 28 | 2 |
| 7 | 18 | 1 | 7 | 18 | 1 |
| 8 | 29 | | 8 | 29 | |
| 9 | 16 | 2 | 9 | 53 | |
| 10 | 39 | | 10 | 39 | |

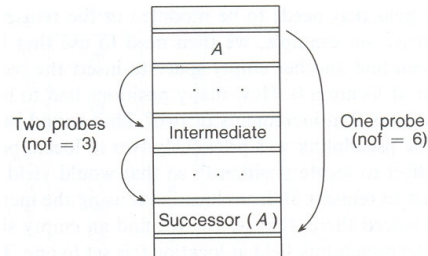
Home address → x
y
z (Immediate predecessor)
Temporary pointer → (Removed item)

Direct File Organization

Computed Chaining

Değerlendirme

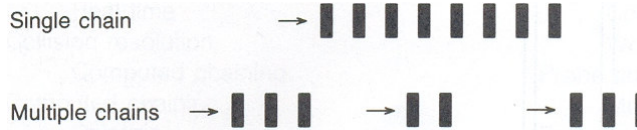
- Link kullanmayan metotlara göre daha iyi performansa sahiptir
- Bir kaydın retrieval probe sayısı probe chain'deki sırasına eşittir
- Worst case retrieval performansı herhangi bir home adresteki çakışma sayısı ile belirlenir
- Worst case retrieval performansı daha düzgün dağılımlı hash fonksiyonu seçimiyle iyileştirilebilir
- Silme işlemi silinen elemandan sonraki tüm kayıtların yeniden yerleştirilmesini gerektirir
- Computed chain metodu diğer metotlara göre silme ve ekleme için daha fazla süre gerektirir ancak retrieval için daha az süre gerektirir



Direct File Organization

Computed Chaining – Multiple Chains

- Divide and Conquer stratejisini kullanır
- Bir probe chain birden fazla parçaya bölünür ve bir kayıt için sadece ilgili olduğu gruba bakılır
- Bir kayıt alanına birden fazla pseudolink oluşturulur
- Standart hash fonksiyonu ve artırma fonksiyonunun yanısıra birde probe chain seçimi için kullanılan $g(\text{key})$ fonksiyonu vardır
- $g(\text{key}) = 0, 1, \dots, R-1$ olabilir. R alt grup sayısıdır
- $g(\text{key}) = \text{key} \bmod R$



Direct File Organization

Computed Chaining – Multiple Chains

Örnek

- 27, 18, 29, 28, 39, 13, 16, 38 ve 53
- $\text{Hash}(\text{key}) = \text{key} \bmod 11$
- $i(\text{key}) = \text{Quotient}(\text{key} / 11) \bmod 11$
- $g(\text{key}) = \text{key} \bmod 2$

| | Key | nof | |
|----|-----|-----|---|
| | | 0 | 1 |
| 0 | 16 | 1 | |
| 1 | 38 | | |
| 2 | 13 | | |
| 3 | | | |
| 4 | | | |
| 5 | 27 | 3 | |
| 6 | 28 | | 2 |
| 7 | 18 | | 1 |
| 8 | 29 | | |
| 9 | 53 | | |
| 10 | 39 | | |

49 eklendi

| | Key | nof | |
|----|-----|-----|---|
| | | 0 | 1 |
| 0 | 16 | 1 | |
| 1 | 38 | | |
| 2 | 13 | | |
| 3 | | | |
| 4 | 49 | | |
| 5 | 27 | 3 | 5 |
| 6 | 28 | | 2 |
| 7 | 18 | | 1 |
| 8 | 29 | | |
| 9 | 53 | | |
| 10 | 39 | | |



Direct File Organization

Comparison of Collision Resolution Methods

Karşılaştırma

COMPARISON OF COLLISION RESOLUTION METHODS

99

TABLE 3.3 COMPARISON OF MEAN NUMBER OF PROBES FOR SUCCESSFUL LOOKUP ($n = 997$; $\alpha =$ packing factor)

| α (percent) | DCWC (10-bit link)* | LISCH (10-bit link) | Progressive overflow† | Linear quotient | Brent's method | Binary tree | Computed chaining (6-bit link) | Computed chaining (2-bit link) |
|-----------------------|---------------------------|---------------------------|--------------------------|--------------------|-------------------|----------------|--------------------------------------|--------------------------------------|
| 20 | 1.100 | 1.106 | 1.125 | 1.15 | 1.102 | 1.102 | 1.070 | 1.070 |
| 40 | 1.200 | 1.232 | 1.333 | 1.277 | 1.217 | 1.217 | 1.168 | 1.214 |
| 60 | 1.300 | 1.379 | 1.750 | 1.527 | 1.367 | 1.364 | 1.264 | 1.381 |
| 70 | 1.350 | — | 2.167 | 1.720 | 1.444 | — | 1.323 | 1.528 |
| 80 | 1.400 | 1.566 | 3.000 | 2.011 | 1.599 | 1.579 | 1.356 | 1.715 |
| 90 | 1.450 | 1.674 | 5.500 | 2.558 | 1.802 | 1.751 | 1.408 | 2.062 |
| 95 | 1.475 | 1.734 | 10.500 | 3.153 | 1.972 | 1.880 | 1.433 | 2.414 |
| 99 | 1.495 | 1.783 | 50.500 | 4.651 | 2.242 | 2.049 | 1.601 | 3.330 |
| 100 | 1.500 | — | — | 6.522 | 2.494 | 2.134 | — | — |

*Theoretical results; mean probes = $1 + \alpha/2$.

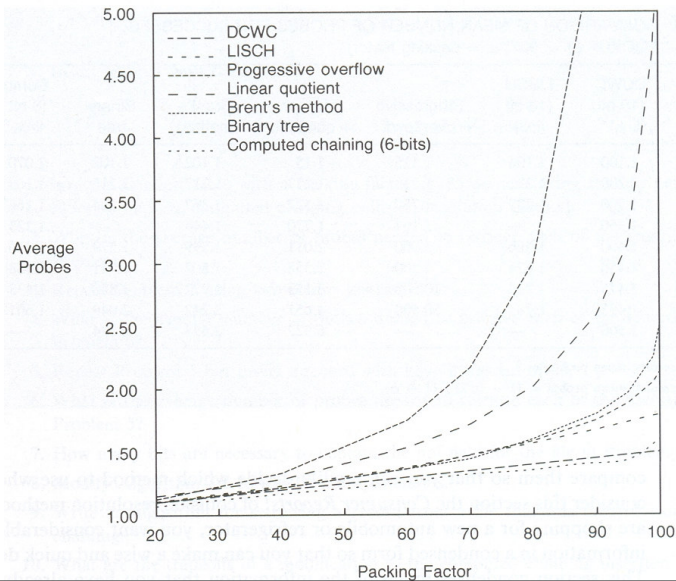
†Theoretical results; mean probes = $(1 - \alpha/2) / (1 - \alpha)$



Direct File Organization

Comparison of Collision Resolution Methods

Karşılaştırma





Direct File Organization

Comparison of Collision Resolution Methods

Karşılaştırma

TABLE 3.4 SEARCH, RELOCATION, AND STORAGE COMPARISONS

| Criteria | LISCH | Progressive overflow | Linear quotient | Brent's method | Binary tree | Computed chaining |
|-------------------|-------|----------------------|-----------------|----------------|-------------|-------------------|
| Successful search | | | | | | |
| Best case | 1 | 1 | 1 | 1 | 1 | 1 |
| Worst case | n | n | n | n | n | n |
| Move an item | No | No | No | Yes | Yes | Yes |
| Extra storage | Yes | No | No | No | No | Yes |



Direct File Organization

Comparison of Collision Resolution Methods

Karşılaştırma

TABLE 3.5 ADVANTAGES, DISADVANTAGES, AND WHEN TO USE VARIOUS COLLISION RESOLUTION METHODS

| Method | Advantages | Disadvantages | When to use |
|--------------------------|---|--|--|
| Coalesced hashing [DCWC] | Excellent performance | Requires space for link field | When ample space is available |
| | Superior performance | Requires space for link field | When ample space is available and insertion time is not critical |
| Progressive overflow | No additional space; simple algorithm | Very poor performance | Almost NEVER |
| Linear quotient | No additional space; reasonable performance | Performance below that of other methods | When space is at a premium; retrieval performance is not critical |
| Brent's method | No additional space; good performance | Performance below that of linked methods | When space is at a premium, even space for a binary tree |
| Binary tree | No additional space; very good performance | Performance below that of linked methods | When file space is limited but temporary space is available for tree; performance is important |
| Computed chaining | Small amount of additional space; excellent performance | Requires some space for pseudolink field; performance below that of DCWC | When some extra file space is available and performance is important; insertion time is not critical |



Direct File Organization

Perfect Hashing

- Perfect hashing bir key değerini unique bir adresle eşleştirir.

hash (key) -> unique address.

- Az sayıdaki kayda sahip dosyalarla kullanılır. (Örn.: Programlama dillerindeki reserved words, en çok kullanılan kelimeler v.b.).
- Eğer potansiyel adres sayısı ile anahtar sayısı aynı ise bu fonksiyona minimum perfect hashing fonksiyonu denir.

$p.\text{hash}(\text{key}) \rightarrow (h_0(\text{key}) + g[h_1(\text{key})] + g[h_2(\text{key})]) \bmod N$

- Worst case computational complexity değeri $O(r^6)$, r toplam kayıt sayısıdır



Direct File Organization

Perfect Hashing – Cichelli's Algorithm

- Aşağıdaki fonksiyonlar kullanılır

$h_0 \rightarrow \text{length}(\text{key})$

$h_1 \rightarrow \text{first_karakter}(\text{key})$

$h_2 \rightarrow \text{last_karakter}(\text{key})$

$g \rightarrow T(x)$

- T herbir x karakteri için hazırlanan tablodaki değeri gösterir
- Ençok zaman alan kısım T değerinin hesaplanmasıdır



Direct File Organization

Perfect Hashing – Cichelli's Algorithm

TABLE 3.6 VALUES ASSOCIATED WITH THE CHARACTERS OF THE PASCAL RESERVED WORDS

| | | | | | |
|--------|--------|--------|--------|-------|--------|
| A = 11 | B = 15 | C = 1 | D = 0 | E = 0 | F = 15 |
| G = 3 | H = 15 | I = 13 | J = 0 | K = 0 | L = 15 |
| M = 15 | N = 13 | O = 0 | P = 15 | Q = 0 | R = 14 |
| S = 6 | T = 6 | U = 14 | V = 10 | W = 6 | X = 0 |
| Y = 13 | Z = 0 | | | | |

$$\begin{aligned} p.\text{hash}(\text{begin}) &= 5 + T(h_1(\text{key})) + T(h_2(\text{key})) \\ &= 5 + T(b) + T(n) \\ &= 5 + 15 + 13 \\ &= 33 \end{aligned}$$



Direct File Organization

Perfect Hashing – Cichelli's Algorithm

- Tablonun oluşturulması için önce bütün anahtarlar ilk ve son karakterlerinin kullanım sıklığı belirlenir
- Herbir anahtara ilk ve son karakterlerinin kullanım sıklıklarının toplam değeri atanır
- Anahtarlar atanan değerlere göre büyükten küçüğe sıralanır
- Backtracking metodu kullanılarak herbir başlangıç ve bitiş karakterine değer atanır

Direct File Organization

Perfect Hashing – Cichelli's Algorithm

Örnek

- cat, ant, dog, gnat, chimp, rat, toad
- Bir karaktere en fazla 4 değeri atanır

İlk ve son karakterlerin sıklığı:

a = 1 c = 2 d = 2 g = 2 p = 1 r = 1 t = 5

| | | | | | | |
|--------------|---|--------------|---|--|--------------|---|
| <u>cat</u> | 7 | <u>toad</u> | 7 | İlk ve son karakterleri daha önce geçenler en son geçenden sonraya alınır | <u>toad</u> | 7 |
| <u>ant</u> | 6 | <u>gnat</u> | 7 | | <u>gnat</u> | 7 |
| <u>dog</u> | 4 | <u>cat</u> | 7 | | <u>dog</u> | 4 |
| <u>gnat</u> | 7 | <u>rat</u> | 6 | | <u>cat</u> | 7 |
| <u>chimp</u> | 3 | <u>ant</u> | 6 | | <u>rat</u> | 6 |
| <u>rat</u> | 6 | <u>dog</u> | 4 | | <u>ant</u> | 6 |
| <u>toad</u> | 7 | <u>chimp</u> | 3 | | <u>chimp</u> | 3 |

Direct File Organization

Perfect Hashing – Cichelli's Algorithm

Örnek (devam)

Bütün anahtarlar sırayla değerlendirilir. Aynı olanlar için back tracking yapılarak ilk ve son karakterlerine yeniden değer atanır.

| | | | | | | | | | |
|------------------|---------|------------------|------------------|------------------|------------------|---------|------------------|---------|---------|
| $t = 0$ | $d = 0$ | $t = 0$ | $d = 0$ | $g = 1$ | $t = 0$ | $d = 0$ | $g = 2$ | | |
| p.hash(toad) = 4 | | p.hash(toad) = 4 | | p.hash(gnat) = 5 | p.hash(toad) = 4 | | p.hash(gnat) = 6 | | |
| $t = 0$ | $d = 0$ | $g = 2$ | $t = 0$ | $d = 0$ | $g = 2$ | $c = 0$ | | | |
| p.hash(toad) = 4 | | | p.hash(toad) = 4 | | p.hash(gnat) = 6 | | | | |
| p.hash(gnat) = 6 | | | p.hash(gnat) = 6 | | p.hash(dog) = 5 | | | | |
| p.hash(dog) = 5 | | | p.hash(dog) = 5 | | p.hash(cat) = 3 | | | | |
| $t = 0$ | $d = 0$ | $g = 2$ | $c = 0$ | $r = 4$ | $t = 0$ | $d = 0$ | $g = 3$ | $c = 0$ | $r = 2$ |
| p.hash(toad) = 4 | | | | | p.hash(toad) = 4 | | | | |
| p.hash(gnat) = 6 | | | | | p.hash(gnat) = 7 | | | | |
| p.hash(dog) = 5 | | | | | p.hash(dog) = 6 | | | | |
| p.hash(cat) = 3 | | | | | p.hash(cat) = 3 | | | | |
| p.hash(rat) = 7 | | | | | p.hash(rat) = 5 | | | | |



Direct File Organization

Perfect Hashing – Cichelli's Algorithm

Örnek (devam)

$t = 0$ $d = 0$ $g = 4$ $c = 0$ $r = 2$ $a = 3$

p.hash(toad) = 4
p.hash(gnat) = 8
p.hash(dog) = 7
p.hash(cat) = 3
p.hash(rat) = 5
p.hash(ant) = 6

Son çözüm aşağıdaki gibidir

$t = 0$ $d = 0$ $g = 4$ $c = 0$ $r = 2$ $a = 3$ $p = 4$

p.hash(toad) = 4
p.hash(gnat) = 8
p.hash(dog) = 7
p.hash(cat) = 3
p.hash(rat) = 5
p.hash(ant) = 6
p.hash(chimp) = 9



Direct File Organization

Perfect Hashing – Cichelli's Algorithm

Değerlendirme

- Toplam backtracking sayısı hesaplama süresini çok etkilemektedir
- Özellikle gerçek zamanlı uygulamalarda minimal perfect function oluşturmak için geçen süre önemsenmez
- Dezavantajları:
 - Aynı uzunlukta ve aynı ilk ve son karaktere sahip iki anahtar olamaz
 - 45 elemandan fazla olan listeleri alt gruplara bölmek gerekebilir
 - Her liste için minimal perfect hashing function bulunamayabilir
- Algoritmanın programlanması kolaydır



Direct File Organization

Haftalık Ödev

C#.NET programlama dilindeki reserved word'ler için Cichelli's Algorithm kullanarak perfect hash fonksiyonu oluřturunuz. Program görsel arayüze sahip olacak ve girilen reserved word için üretilen adres deęerini ekrana yazacak. Ayrıca ayrılmıř kelimelerin oluřturacađınız dosya içindeki sıralanıřları ekranda görüntülenecek.