

Dictionaries

Datastrukturer (recap)

Datastruktur = data + operationer herpå

Data:

- ▶ En key (nøgle) + associeret data (ofte underforstået, også i disse slides).

Operationer:

- ▶ Datastrukturens egenskaber udgøres af **de tilbudte operationer** (API for adgang til data), samt **deres køretider** (forskellige implementationer af samme API kan give forskellige køretider).

I dette kursus: katalog af **datastrukturer (niveau 1)** med bred anvendelse samt **effektive implementationer (niveau 2)** heraf.

Datastrukturer (recap)

Vi har allerede set **Priority Queue**. Datastruktur som tilbyder følgende operationer:

- ▶ **Extract-Min()**: Fjern et element med mindste nøgle fra prioritetskøen og returner det.
- ▶ **Insert(key)**: Tilføj nyt element til prioritetskøen.
- ▶ **Build(liste af elementer)**: Byg en prioritetskø indeholdende elementerne.
- ▶ **Decrease-Key(key, reference til element i kø)**: Sætter nøglen for elementet til $\min\{\text{key}, \text{gamle key}\}$.

Dictionaries

Emne for disse slides: **Dictionary**. Datastruktur som tilbyder følgende operationer:

- ▶ **Search(key)**: returner element med nøglen key (eller fortæl hvis det ikke findes).
- ▶ **Insert(key)**: Indsæt nyt element med nøglen key.
- ▶ **Delete(key)**: Fjern element med nøglen key.

- ▶ **Predecessor(key)**: Find elementet med højeste nøgle $< \text{key}$.
- ▶ **Successor(key)**: Find elementet med laveste nøgle $> \text{key}$.
- ▶ **OrderedTraversal()**: Udskriv elementer i sorteret orden.

For de sidste tre operationer bruges nøglernes ordning.

Hvis kun de tre første operationer skal understøttes, kaldes det en **unordered dictionary**. Hvis alle seks understøttes, kaldes det en **ordered dictionary**.

Vi antager i disse slides unikke nøgler (ingen dubletter blandt elementer).

Dictionaries

Dictionaries i C#: interface IDictionary.

Dictionaries i Java: interface Map.

Dictionaries i Python: dict.

Implementationer (niveau 2) som vi møder i dette kursus:

- ▶ **Balancerede binære søgetræer**: Understøtter alle ovenstående operationer (samt mange flere, f.eks. ved at tilføje ekstra information i knuderne) i $O(\log n)$ tid.
- ▶ **Hashing**: understøtter de tre første operationer i $O(1)$ tid (medmindre man er uheldig).

Disse implementationer findes i C# som SortedDictionary og Dictionary. De findes i Java som TreeMap og HashMap. I Python er den indbyggede datatype dict implementeret med hashing (der er ingen balancerede binære søgetræer i Pythons standard moduler, men man kan finde moduler med dem fra andre kilder).

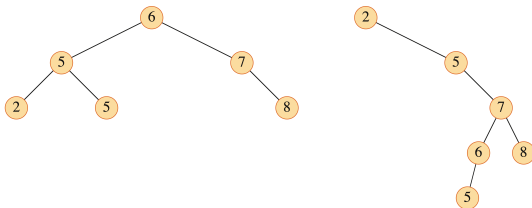
Binært søgetræ

- ▶ et binært træ
- ▶ med knuder i *inorder*

Et binært træ med nøgler i alle knuder overholder *inorder* hvis det for alle knuder v gælder:

nøgler i v 's *venstre undertræ* \leq nøgle i v \leq nøgler i v 's *højre undertræ*

Eksempler:

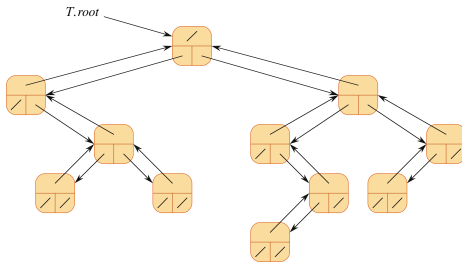


Binære søgetræer

Typisk implementation: **Knude-objekter** med:

- ▶ Reference til forælder
- ▶ Reference til venstre undertræ
- ▶ Reference til højre undertræ

samt ét **træ-objekt** med reference til roden. (Programmeringsprog: reference, bog: pointer).



Knude-objekter

Java-klasse for knuder:

```
class Node {  
    int key;  
    Node leftchild;  
    Node rightchild;  
    Node parent;  
    .  
    .  
    (constructor)  
    (andre metoder)  
    .  
}
```

Python-klasse for knuder:

```
class Node:  
    def __init__(self):  
        self.key = None  
        self.leftchild = None  
        self.rightchild = None  
        self.parent = None  
    .  
    .  
    (andre metoder)  
    .
```

C#-klasse for knuder: Ligner Java-eksemplet.

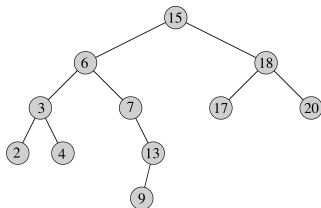
Binære søgetræer

Pga. definitionen af inorder

nøgler i v 's venstre undertræ \leq nøgle i $v \leq$ nøgler i v 's højre undertræ

kan binære søgetræer siges at indeholde data i sorteret orden.

Mere præcist: **inorder gennemløb** vil udskrive nøgler i sorteret orden:



INORDER-TREE-WALK(x)

if $x \neq \text{NIL}$

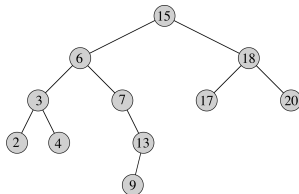
INORDER-TREE-WALK($x.\text{left}$)

print $\text{key}[x]$

INORDER-TREE-WALK($x.\text{right}$)

Køretid: $O(n)$ (der laves $O(1)$ arbejde per knude i træet).

Søgning i binære søgetræer



TREE-SEARCH(x, k)

if $x == \text{NIL}$ or $k == \text{key}[x]$

return x

if $k < x.\text{key}$

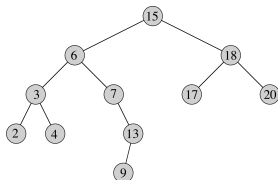
return TREE-SEARCH($x.\text{left}, k$)

else return TREE-SEARCH($x.\text{right}, k$)

Princip (\Rightarrow korrekthed):

Hvis søgte element findes, er det i det undertræ, vi er kommet til

Flere slags søgninger i binære søgetræer



TREE-MAXIMUM(x)

```
while  $x.right \neq \text{NIL}$   
     $x = x.right$   
return  $x$ 
```

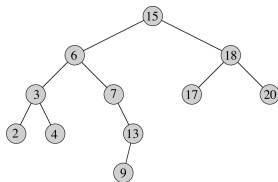
TREE-MINIMUM(x)

```
while  $x.left \neq \text{NIL}$   
     $x = x.left$   
return  $x$ 
```

Princip (\Rightarrow korrekthed):

Det søgte element findes i det undertræ, vi er kommet til

Flere slags søgninger i binære søgetræer



TREE-SUCCESSOR(x)

if $x.right \neq \text{NIL}$

return TREE-MINIMUM($x.right$)

$y = x.p$

while $y \neq \text{NIL}$ and $x == y.right$

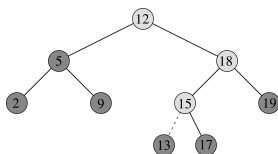
$x = y$

$y = y.p$

return y

Korrekthed: Via inorder kan man udelukke de fleste knuder i træet fra at være x 's successor. Tilbage bliver mindste knude i x 's højre undertræ, og ellers første forfader, som har x i venstre undertræ.

Indsættelser i binære søgetræer



- ▶ Søg nedad fra rod: gå i hver knude v mødt videre ned i det undertræ (højre/venstre), hvor nye element skal være iflg. inorder-krav for v .
- ▶ Når blad (NIL/tomt undertræ) nås, erstat dette med den nye knude (med to tomme undertræer).

Inorder er overholdt for knuder på søgesti (pga. søgeregul), og for alle andre knuder (fordi de ikke har fået nogle nye efterkommere i deres to undertræer).

Sletninger i binære søgetræ

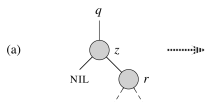
Sletning af knude z :

- ▶ Case 1: Mindst ét barn er NIL: Fjern z samt dette barn, lad andet barn tage z 's plads.
- ▶ Case 2: Ingen børn er NIL: Da er successor-knuden y til z den mindste knude i z 's højre undertræ. Fjern y (som er en Case 1 fjernelse, da dens venstre barn er NIL), og indsæt den på z 's plads.

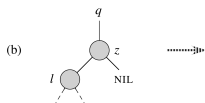
Begge cases efterlader træet i inorder: I Case 1 får ingen knuder nye efterkommere i deres to undertræer. I Case 2 får y (og kun y) nye efterkommere i sine to undertræer, men da y er z 's successor, er der ingen nøgler i træet med værdi mellem z 's og y 's nøgler, så nøglerne i y 's nye undertræer overholder inorder i forhold til y , eftersom de gjorde i forhold til z .

Bemærk at strukturelt i træet er alle sletninger en Case 1 sletning.

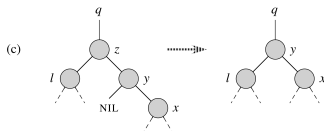
Sletninger i binære søgetræ (bogens cases)



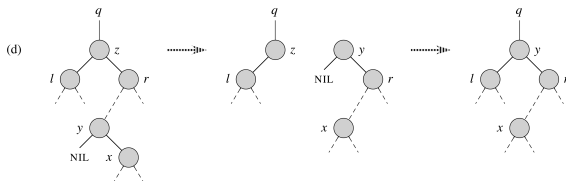
Case 1



Case 1



(Case 2 \rightarrow) Case 1



(Case 2 \rightarrow) Case 1

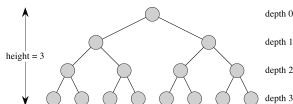
Tid for operationer i binære søgetræ

For alle operationer (undtagen inorder gennemløb):

Gennemløb sti fra rod til blad.

Dvs. køretid = $O(\text{højde})$.

Et træ med højde h kan ikke indeholde flere knuder end det fulde træ med højde h . Dette indeholder $2^{h+1}-1$ knuder (jvf. slides om heaps).



Så for et træ med n knuder og højde h gælder:

$$n \leq 2^{h+1} - 1 \quad \Leftrightarrow \quad \log_2(n + 1) - 1 \leq h$$

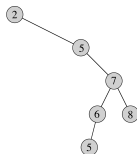
Dvs. den bedst mulige højde er $\log_2 n (\pm 1)$

Kan vi holde højden tæt på optimal – f.eks. $O(\log n)$ – under updates (indsættelser og sletninger)?

Balancerede binære søgetræer

Kan vi holde højden $O(\log n)$ under updates (indsættelser og sletninger)?

Kræver **rebalancing** (omstrukturering af træet) efter updates, da dybe træer ellers kan opstå:



Vedligehold af $O(\log n)$ højde første gang opnået med AVL-træer [1961].

Mange senere forslag. Et forslag består af:

- ▶ Balanceinformation i knuder.
- ▶ Strukturkrav baseret på balanceinformation, som sikrer $O(\log n)$ højde.
- ▶ Algoritmer, som genopretter strukturen efter en update.

I dette kursus: **rød-sortede træer**.

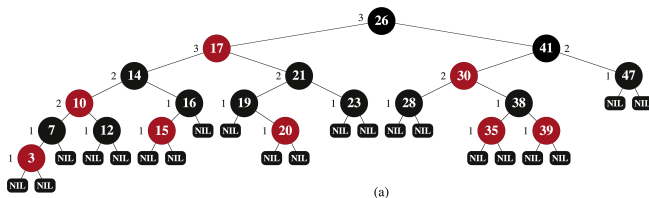
Rødsorte træer

Balanceinformation i knuder: 1 bit (kaldet rød/sort farve).

Strukturkrav:

- ▶ Rod og blade sorte.
- ▶ Samme antal sorte på alle rod-blad stier.
- ▶ Ikke to røde i træk på nogen rod-blad sti.

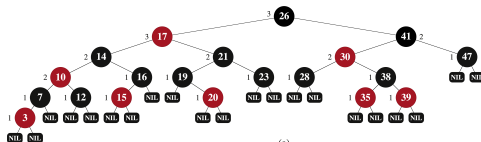
Eksempel:



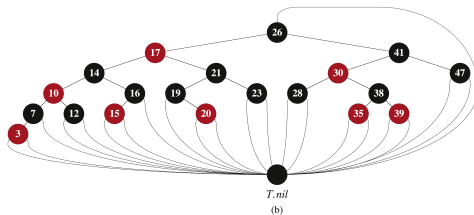
NB: begrebet blade bliver i rød-sorte træer af brugt om NIL-undertræer (hvilket teknisk set øger højden af et træ med én).

Rødsorte træer

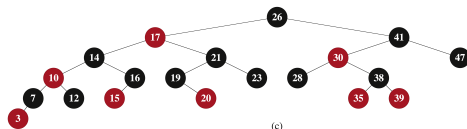
Andre repræsentationer i bogen (samme træ):



(a)



(b)



(c)

Rødsorte træer

Strukturkrav (recap):

- ▶ Rod og blade sorte.
- ▶ Samme antal sorte på alle rod-blad stier.
- ▶ Ikke to røde i træk på nogen rod-blad sti.

Sikrer disse strukturkrav sikrer $O(\log n)$ højde? [Ja](#):

Lad antal sorte på alle rod-blad stier være k (dvs. $k - 1$ sorte knuder plus et sort blad). Så indeholder alle rod-blad stier mindst $k - 1$ knuder (de sorte plus evt. nogle røde). Der er derfor mindst $k - 1$ fulde lag af knuder i træet.

Derfor er $n \geq 2^0 + 2^1 + 2^2 + \dots + 2^{k-2} = 2^{k-1} - 1$.

Heraf følger $\log(n + 1) \geq k - 1$.

Da der ikke er to røde knuder i træk, indholder den længste rod-blad sti højst $2(k - 1)$ kanter.

Så højden er højst $2(k - 1) \leq 2 \log(n + 1) = O(\log n)$.

Indsættelse

1. Indsæt en knude i træet
2. Fjern evt. opstået ubalance (overtrædelse af rød-sort strukturkravene).

Recall indsættelse: et blad (NIL) erstattes af en knude med to blade som børn.

Ubalance?

Recall indsættelse: et blad (NIL) erstattes af en knude med to blade som børn.

Overtrædelse af rød-sorter strukturkrav?

- ▶ Rod og blade sorte.
- ▶ Samme antal sorte på alle rod-blad stier.
- ▶ Ikke to røde i træk på nogen rod-blad sti.

De to nye blade skal være sorte.

Vi vælger at gøre den nye indsatte knude rød.

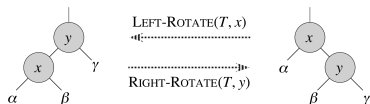
Mulige overtrædelse af strukturkrav er nu: To røde knuder i træk på en rod-blad sti ét sted i træet.

Idé til plan: Kan problemet ikke løses umiddelbart, så skub det opad i træet indtil det kan (forhåbentligt nemt at gøre, hvis det når roden).

Rebalancing

Detaljeret plan: skub rød-rød problem opad i træet, under brug af omfarvninger og restruktureringer af træet.

Den basale restrukturering vil være en **rotation** (α, β, γ er undertræer, evt. tomme):



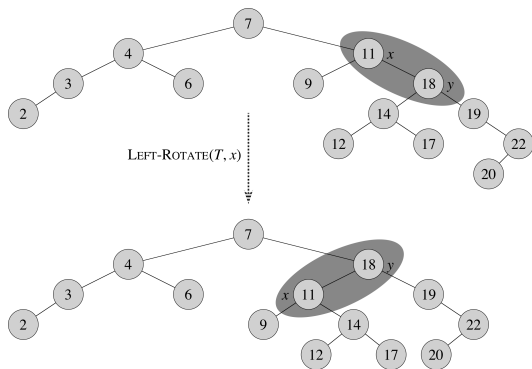
Central observation: Rotationer kan ikke ødelægge in-order i træet:

Kun x og y kan få in-order overtrådt (alle andre knuder har de samme elementer i deres undertræer før og efter). Men dette sker ikke, da følgende gælder både før og efter en rotation:

$$\text{keys i } \alpha \leq x \leq \text{keys i } \beta \leq y \leq \text{keys i } \gamma$$

Så vi skal ikke bekymre os om bevarelse af in-order, hvis vi kun restrukturerer vha. rotationer.

Eksempel på rotation



Plan for rebalancering (efter indsættelse)

Recap af plan: skub rød-rød problem opad i træet, under brug af omfarvninger og restruktureringer (rotationer) af træet.

Princip undervejs:

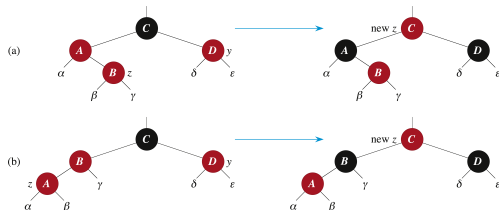
- ▶ To røde knuder i træk på en rod-blad sti højst ét sted i træet.
- ▶ Bortset herfra er de rød-sortे krav overholdt.

Mål: I $O(1)$ tid, fjern problemet eller skub det ét skridt nærmere roden.

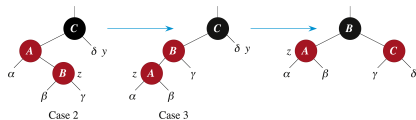
Dette vil give rebalancering i $O(\text{højde}) = O(\log n)$ tid.

Cases i rebalancing (efter indsættelse)

Case I: Rød onkel til z (onkel = forælders søsken y).



Case II: Sort onkel til z (onkel = forælders søsken y).

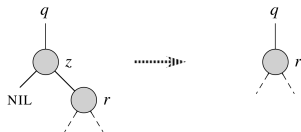


Her er z nederste knude i rød-rød problemet. Kontrollér at princip vedligeholdes. Kontrollér at problem fjernes eller flyttes nærmere roden (én færre sort knude på sti til rod). Hvis z bliver lig roden, kan den blot farves sort (\Rightarrow alle stier får en sort mere).

Sletning

1. Slet en knude i træet
2. Fjern evt. opstået ubalance (overtrædelse af rød-sort kravene).

Recall sletning: der fjernes strukturelt set altid én knude hvis ene barn er et blad (NIL), som også fjernes.



Ubalance?

Overtrædelse af rød-sorter krav?

- ▶ Rod og blade sorte.
- ▶ Samme antal sorte på alle rod-blad stier.
- ▶ Ikke to røde i træk på nogen rod-blad sti.

Fjernet knude rød: Alle rød-sorter krav stadig overholdt.

Fjernet knude sort: Ikke længere samme antal sorte på alle stier.

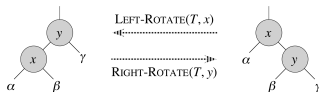
Meget brugbar formulering:

Lad den fjernede knudes andet barn være "sværtet" og gælde for "én mere" sort end dens farve angiver, når vi tæller sorte på stier (sværtet sort = 2 sorte, sværtet rød = 1 sort). Så er kravene overholdt, bortset fra eksistensen af en sværtet knude.

Idé til plan: Kan problemet ikke løses umiddelbart, så skub det opad i træet til det kan (forhåbentligt nemt at gøre, hvis det når roden).

Rebalancing

Skub sværtet knude opad i træet, under brug af omfarvninger og rotationer:



Princip undervejs:

- ▶ Højst én knude i træet er sværtet.
- ▶ Hvis sværtningen tælles med, er de rød-sortे krav overholdt.

Nemme stoptilfælde:

- ▶ Sværtet knude er rød \Rightarrow sværtning kan fjernes ved at farve knuden sort.
- ▶ Sværtet knude er rod \Rightarrow sværtning kan bare fjernes (\Rightarrow alle stier får en sort mindre).

(Så nok at se på tilfældet at den sværtede knude er sort.)

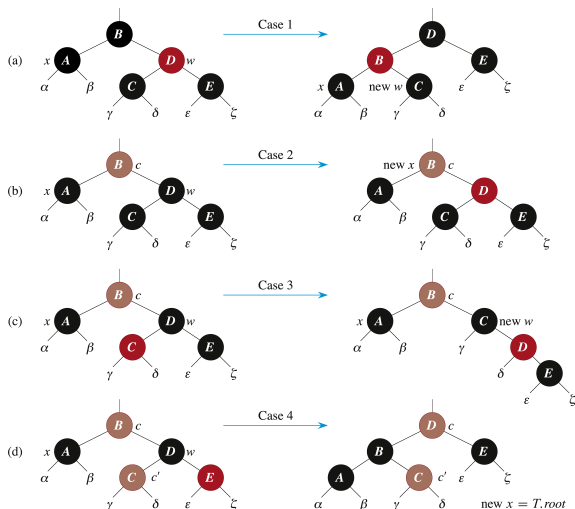
Rebalancing

Mål: I $O(1)$ tid, fjern problemet eller skub det ét skridt nærmere roden.
Dette vil give rebalancing i $O(\text{højde}) = O(\log n)$ tid.

Cases for sværtet sort knude x med søsken w .

1. Rød søsken.
2. Sort søsken, og denne har to sorte børn.
3. Sort søsken, og dennes nærmeste barn er rødt, det fjerneste sort.
4. Sort søsken, og dennes fjerneste barn er rødt.

Cases i rebalancing (efter sletning)



Her er x svættet knude. Kontrollér at princip vedligeholdes. Kontrollér $O(1)$ tid før svættning fjernes eller flyttes ét skridt nærmere roden.

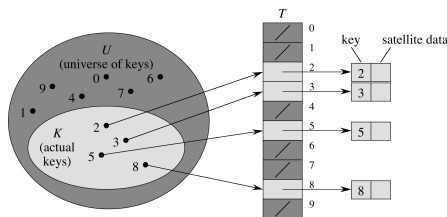
Hashing

Vi antager i hashing, at keys er **heltal** op til en max-grænse u . [For at bruge hashing på andre datatyper må elementer tildeles en unik heltalsværdi, jvf. `GetHashCode()` i C#, `hashCode()` i Java og `hash()` i Python.]

Dvs. at vi har et univers af mulige heltalskeys: $U = \{0, 1, 2, \dots, u - 1\}$.

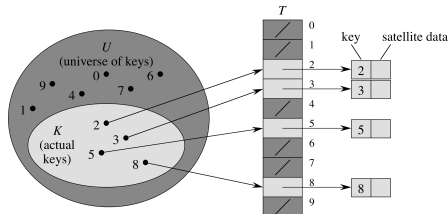
En dictionary gemmer en delmængde $K \subseteq U$, f.eks. $K = \{2, 3, 5, 8\}$

Udgangspunktet i hashing er (lige som i Counting sort) idéen om at bruge keys som indekser i et array:



Dette tager $O(1)$ tid for Search, Insert og Delete.

Problem: pladsforbrug



Problem: Denne idé kan nemt generere pladsspild, fordi array-størrelse er $u = |U|$, mens antallet $n = |K|$ af gemte elementer ofte er meget mindre.

Eksempel: gem 5 CPR-numre. Dvs. keys af typen 180781-2345, der kan ses som heltal i $U = \{0, 1, 2, \dots, 10^{10} - 1\}$. Her er $u = 10^{10}$, mens $n = 5$. Så vi bruger mindst 10^{10} bytes (> 8 Gb RAM) for at gemme 5 CPR-numre.

Gemmes 32 eller 64 bits heltal er $u = 2^{32} \approx 10^{10}$ eller $u = 2^{64} \approx 10^{20}$. Dvs. samme situation eller meget værre.

Hash-funktioner

En løsning på problemet med pladsforbruget: find en funktion h , som mapper keys fra det store heltals-univers U til et mindre:

$$h : U \rightarrow \{0, 1, 2, \dots, m - 1\}.$$

Her er m den ønskede array-størrelse. Ofte vælges $m = O(n)$, så bruges der ikke mere plads på array'et end på elementerne selv.

En sådan funktion kaldes en **hash-funktion**. Et eksempel på en hash-funktion kan være:

$$h(k) = k \bmod m$$

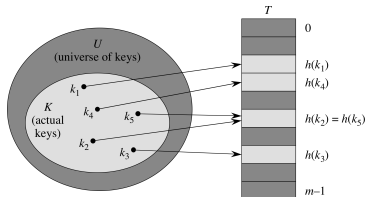
Konkret eksempel med $m = 41$:

$$\begin{array}{ll} h(12) = 12 \bmod 41 = 12 & (\text{da } 0 \cdot 41 + 12 = 12) \\ h(100) = 100 \bmod 41 = 18 & (\text{da } 2 \cdot 41 + 18 = 100) \\ h(479869) = 479869 \bmod 41 = 5 & (\text{da } 11704 \cdot 41 + 5 = 479869) \end{array}$$

Kollisioner

Hashfunktioner løser problemet med pladsforbrug. Men de genererer et andet problem: To keys kan hash'er til samme array index.

$$h(479869) = 479869 \bmod 41 = 5 \quad (\text{da } 11704 \cdot 41 + 5 = 479869)$$
$$h(46) = 46 \bmod 41 = 5 \quad (\text{da } 1 \cdot 41 + 5 = 46)$$

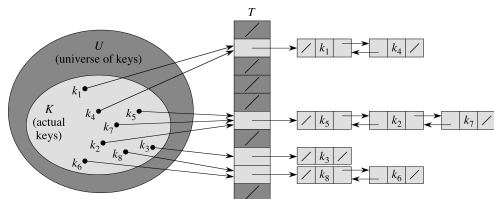


Dette kaldes en **kollision**.

Når h skal mappe U ind i en mindre mængde $\{0, 1, 2, \dots, m-1\}$, vil der altid være nogle keys k og k' hvor $h(k) = h(k')$. Så kollisioner kan ikke undgås, og vi skal derfor finde en løsning.

Chaining

En simpel løsning: en array-indgang indeholder starten på en lænket liste med alle de indsatte elementer, hvis key hash'er til denne array-indgang. Det kaldes **chaining**.



Prisen er, at lænkede lister skal gennemløbes under Search og Delete, hvorved tiden for disse operationer stiger fra $\Theta(1)$ til $\Theta(|\text{liste}|)$. Insert er stadig $O(1)$, da vi bare kan indsætte forrest i listen.

Open addressing

Open addressing er et alternativ til chaining: Forsøg at finde tom slot i selve array'et.

0	
1	79
2	
3	
4	69
5	98
6	
7	72
8	
9	14
10	
11	50
12	

Linear probing (alias linear hashing):

$$h(k, i) = (h_1(k) + i) \bmod m$$

Quadratic hashing:

$$h(k, i) = (h_1(k) + c_1 \cdot i + c_2 \cdot i^2) \bmod m$$

[Ikke med i Cormen et al. 4. udgave, udgår af pensum.]

Double hashing

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

Her er $h_1(k)$ og $h_2(k)$ to hash-funktioner (kaldet “auxiliary” i bog).

Insert: $i = 0, 1, 2, \dots$ forsøges til en empty slot findes.

Search: $i = 0, 1, 2, \dots$ forsøges til element eller empty slot findes.

Open addressing, bemærkninger

0	
1	79
2	
3	
4	69
5	98
6	
7	72
8	
9	14
10	
11	50
12	

Linear probing (alias linear hashing):

$$h(k, i) = (h_1(k) + i) \bmod m$$

Double hashing

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

- ▶ At implementere Delete er mere kompliceret. En simpel løsning: lad slettede elementer stå, mærk dem som slettede, ryd op en gang i mellem ved at genbygge hashtabellen (indsæt de tilbageværende elementer forfra). For linear hashing findes en mere direkte metode, se Cormen et al. 4. udgave, sektion 11.5.1 (ikke pensum).
- ▶ Det er nødvendigt at $n \leq m$ (da alle n elementer ligger i array'et). Gerne $n \approx m/4$ for at undgå at køretider degenererer.
- ▶ De gennemsøgte indexer $\{h(k, 0), h(k, 1), h(k, 2), \dots, h(k, m-1)\}$ bør være $\{0, 1, 2, \dots, m-1\}$ for alle k (så hele array'et gennemses). Dette kan f.eks. sikres ved at vælge m som et primtal.

Køretid for hashing

Vi ønsker en hash-funktion h som spreder keys fra det konkrete input godt ud over $\{0, 1, 2, \dots, m-1\}$, så der bliver få kollisioner. Man tænker ofte på gode hash-funktioner, som nogle der mapper keys til indekser på en tilsyneladende tilfældig måde.

Man kan dog for en given hashfunktion altid finde mindst $|U|/m$ (dvs. u/m) keys, som hash'er til samme array-index. Ofte er $u/m > n$, hvad der gør worst case tiden til $\Theta(n)$.

I praksis håber man ofte bare på, at dette ikke sker for ens konkrete input og konkrete valg af hashfunktion. Der findes faktisk metoder til at garantere, at dette sjældent sker - se næste slide om universal hashing.

I praksis kan man regne med at hashing understøtter Search, Insert og Delete i $O(1)$ køretid, medmindre man er meget uheldig.

Man kan desuden sænke worst case tiden til $O(\log n)$ ved at bruge balancerede søgetræer i stedet for lænkede lister i chaining. Dette gør Java, når den lænkede liste bliver stor.

Universal hashing (ikke pensum)

Betragt følgende familie af hashfunktion:

$$h(k) = ((a \cdot k + b) \bmod p) \bmod m,$$

hvor p er et fast primtal større end $|U|$ og a, b er faste, men tilfældigt valgte heltal med $1 \leq a \leq p - 1$ og $0 \leq b \leq p - 1$.

Man kan bevise (ikke pensum) at ovenstående hashfunktion er god i følgende forstand (kaldet universal hashing):

For ethvert datasæt kan vi for et tilfældigt valg af a og b forvente så få kollisioner, at Search, Insert og Delete tager $O(1)$ tid.

I Cormen et al. 4. udgave beskrives en anden familie af hashfunktioner med samme egenskab (formel (11.2) og formel (11.5)). Udregningen af den bruger færre CPU-skridt og kan være et godt valg, hvis man selv skal implementere hashtabeller. Eller se Mikkel Thorups noter om hashing.

Andre anvendelser af hashfunktioner (ikke pensum)

Hashfunktioner har mange anvendelsesområder, udover som implementation (niveau 2) af unordered dictionaries.

Nogle eksempler (ikke pensum):

- ▶ Fingerprint af filer og dokumenter.
- ▶ Digital signatures.
- ▶ Load balancing
- ▶ Coordinated sampling

De forskellige anvendelser har hver deres specifikke krav til hashfunktionerne.