

# Dictionaries

# Datastrukturer (recap)

Datastruktur = data + operationer herpå
---

## Data:

- ▶ En ID (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).

DM507: katalog af **datastrukturer med bred anvendelse** samt **effektive implementationer heraf**.

# Datastrukturer (recap)

Vi har allerede set **Priority queue**. Datastruktur som understøtter operationerne:

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

# Dictionaries

I dag: datastrukturen [Dictionary](#). I Java: interface Map.

Datastruktur som understøtter operationerne:

- ▶ `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  $\leq$  `key`.
- ▶ `Successor(key)`: Find elementet med laveste nøgle  $\geq$  `key`.
- ▶ `OrderedTraversal()`: Udskriv elementer i sorteret orden.

For de sidste tre operationer kræves at nøglerne har en ordning.

Hvis kun de tre første operationer skal understøttes, kaldes det en [unordered dictionary](#). Hvis alle seks understøttes, kaldes det en [ordered dictionary](#).

# Dictionaries

Implementationer vi møder i DM507:

- ▶ **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 forventet tid  $O(1)$ .

Disse implementationer findes i Java som henholdsvis TreeMap og HashMap.

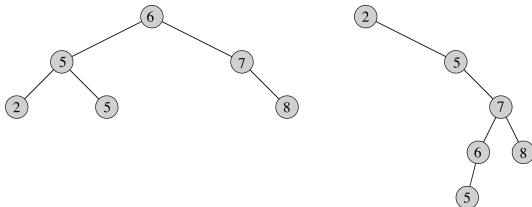
# Binært søgetræ

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

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

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

Eksempler:

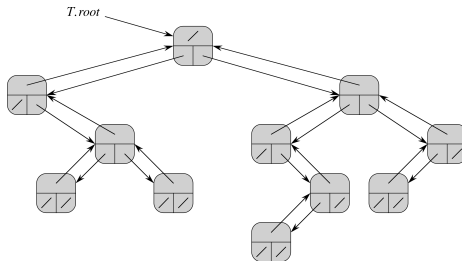


# Binære søgetræer

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

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

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



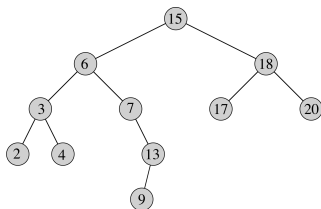
# Binære søgetræer

Pga. definitionen af inorder

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

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( $x.\text{left}$ )

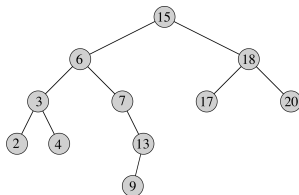
print  $\text{key}[x]$

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

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



# Søgning i binære søgetræer

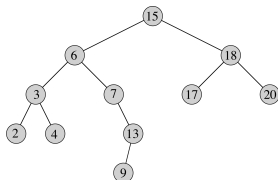


```
TREE-SEARCH( $x, k$ )  
  if  $x == \text{NIL}$  or  $k == \text{key}[x]$   
    return  $x$   
  if  $k < x.\text{key}$   
    return TREE-SEARCH( $x.\text{left}, k$ )  
  else return TREE-SEARCH( $x.\text{right}, k$ )
```

Princip:

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  $x$ 
```

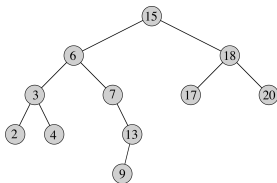
TREE-MINIMUM( $x$ )

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

Princip:

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  $x == y.right$

$x = y$

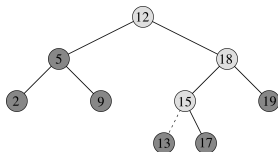
$y = y.p$

**return**  $y$

Princip:

Se på stien fra  $x$  til rod. Ingen side-træer på den kan indeholde det søgte element (pga. in-order).

## 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. søgeregel), og for alle andre knuder (fordi de ikke har fået ændret deres undertræer).

# Sletninger i binære søgetræ

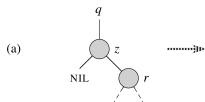
Sletning af knude  $z$ :

- ▶ 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 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.

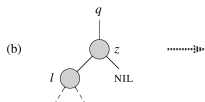
Begge cases efterlader træet i inorder (i Case 2 fordi  $y$  vil overholde inorder når den sættes på  $z$ 's plads, da ingen knuder i træer har nøgle med værdi mellem  $z$ 's og  $y$ 's nøgler).

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

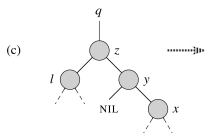
# Sletninger i binære søgetræ (bogens cases)



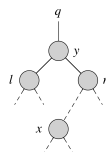
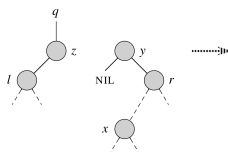
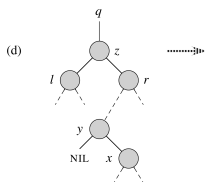
Case 1



Case 1



(Case 2  $\rightarrow$ ) Case 1



(Case 2  $\rightarrow$ ) Case 1

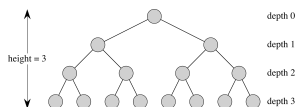
# Tid for operationer i binære søgetræ

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:

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

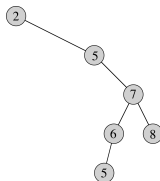
Dvs. den bedst mulige højde er  $\log_2 n (\pm 1)$

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:

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

I DM507: **rød-sortede træer**.



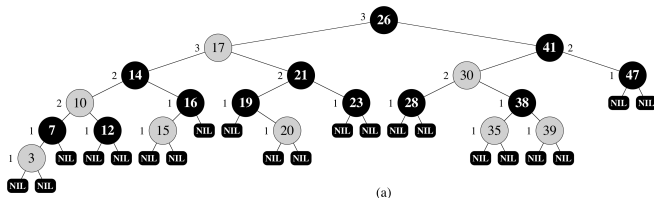
# Rødsorte træer

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.

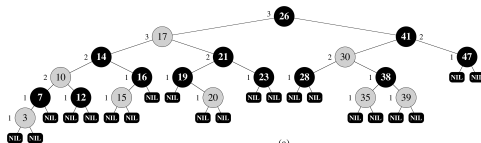
Eksempel:



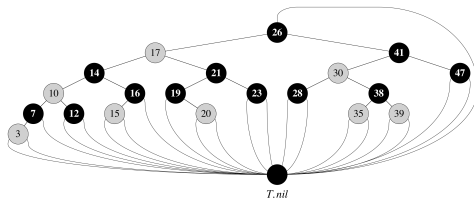
(NB: begrebet blade bliver fra nu af brugt om NIL-undertræer).

# Rødsorte træer

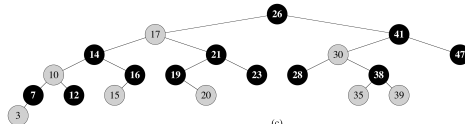
Andre repræsentationer i bogen (samme træ):



(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:

Hvis antal sorte på alle stier er  $k$ , indeholder alle rod-blad stier mindst  $k - 1$  kanter, og der er derfor mindst  $k - 1$  fulde lag af indre knuder.

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

Heraf følger  $\log(n + 1) \geq k$ .

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

Så højde  $\leq 2(k - 1) = 2k - 2 \leq 2\log(n + 1) - 2$ .

# 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 lave nye indsatte knude rød.

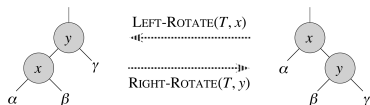
*Mulige overtrædelse af strukturkrav er nu: To røde knuder i træk på en rod-blad sti ét sted i træet.*

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

# Rebalancing

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

Den basale restrukturering vil være en **rotation** ( $\alpha, \beta, \gamma$  er undertræer, evt. tomme):



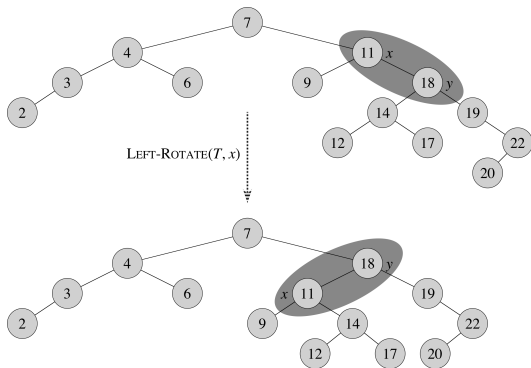
Central observation: Rotationer kan ikke ødelægge in-order i træet:

Kun  $x$  og  $y$  kan få in-order overtrådt (alle andre undertræer indeholder de samme elementer), 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-sort krav overholdt.

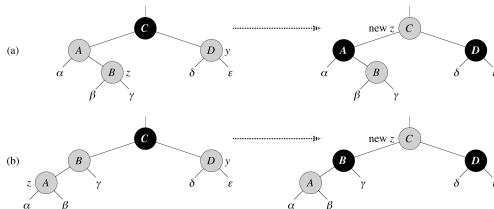
Mål: I  $O(1)$  tid, fjern problemet eller skub det ét skridt nærmere roden.

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

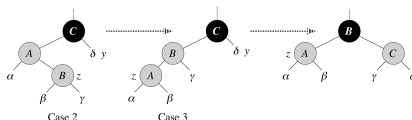


# Cases i rebalancing (efter indsættelse)

Case 1: Rød onkel (onkel = forælders søskend).



Case 2: Sort onkel (onkel = forælders søskend).

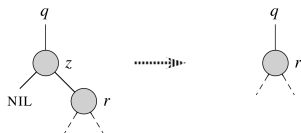


Her er  $z$  nederste knude i rød-rød problemet. Kontrollér at princip vedligeholdes. Kontrollér  $O(1)$  tid før problem fjernes eller flyttes nærmere roden. Når  $p.z$  (eller  $z$ ) er lig roden, kan denne blot farves sort ( $\Rightarrow$  alle stier får en sort mere).

# 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.

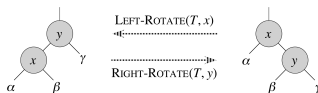
Meget brugbar formulering:

*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.)

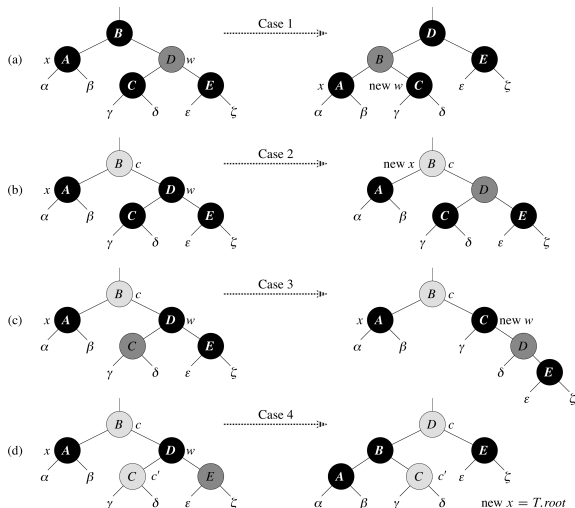
# Rebalancing

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 søskend  $w$ .

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

# Cases i rebalancing (efter sletning)



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

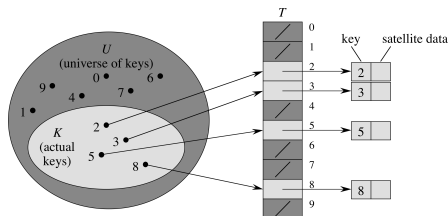
# Idé i hashing

Vi antager keys er heltal (for andre datatyper må man tildele dem en heltalsværdi, jvf. hashCode() i java) op til en max-grænse.

Universe  $U = \{0, 1, 2, \dots, |U| - 1\}$

Dictionary gemmer delmængde  $K \subseteq U$ , e.g.  $K = \{2, 3, 5, 8\}$

Grundideen i hashing er at se på keys som indekser i et array:



Problem: Pladsspild fordi tabelstørrelse er lig  $|U|$ , mens  $n = |K|$  ofte er meget mindre. Gemmes f.eks. 100 CPR-numre er  $|U| = 10^{10}$  mens  $n = 100$ . Gemmes ints er  $|U| = 2^{32}$ .

# Hash-funktioner

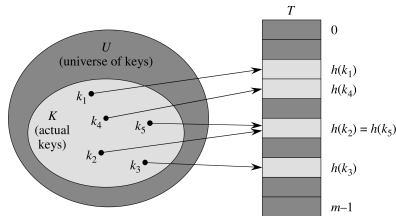
Hash-funktion:

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

Her er  $m$  den ønskede tabel størrelse. Ofte vælges  $m = O(n)$ .

F.eks.:

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





# Hash-funktioner

Endnu bedre:

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

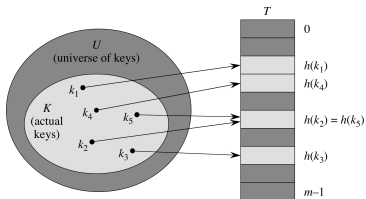
hvor  $p$  er et fast primtal større end  $|U|$  og  $a, b$  er faste, men tilfældigt valgte heltal med  $1 \leq a \leq p$  og  $0 \leq b \leq p$ .

Analyse heraf senere i studiet viser at kvalitet kan garanteres i en bestemt forstand (kaldet universal hashing). Lidt forsimplet genfortalt: vi kan forvente så få kollisioner (se næste side) at alle operationer (Search, Insert, Delete) tager  $O(1)$  tid.

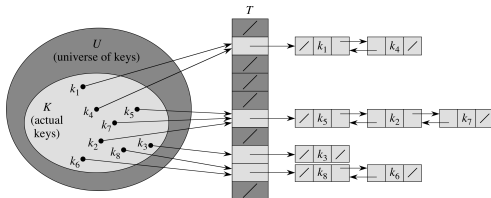
I DM507-bog (inden for DM507 pensum): flere forslag til hash-funktioner, baseret på "erfaring", men alle uden teoretisk garanti for kvalitet.

# Kollisioner

To keys hash'er til samme array index:



En simpel løsning: Chaining (lænkede lister).



# Open addressing

En anden løsning: Forsøg at finde tom slot.

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

Linear hashing:

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

Quadratic hashing:

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

Double hashing

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

Her er  $h'(k)$  og  $h''(k)$  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.

Sletninger: besværligt (simpleste løsning: lad slettede elementer stå, mærk dem som slettede, ryd op en gang i mellem ved at genbygge tabel).

- ▶ Det er nødvendigt at  $n \leq m$  (da alle  $n$  elementer ligger i tabellen). Gerne  $n < m/4$  for at undgå at køretider degenererer.
- ▶  $\{h(k, 0), h(k, 1), h(k, 2), \dots, h(k, m-1)\}$  bør være  $\{0, 1, 2, \dots, m-1\}$  for alle  $k$  (så hele tabellen gennemses).