

Merging og hashing

Mål

Målet for disse slides er at diskutere nogle metoder til at gemme og hente data effektivt.

Dette emne er et uddrag af kurset *DM507 Algoritmer og datastrukturer* (2. semester).

Mål

Målet for disse slides er at diskutere nogle metoder til at gemme og hente data effektivt.

Dette emne er et uddrag af kurset *DM507 Algoritmer og* *datastrukturer* (2. semester).

Datastrukturer

Datastruktur = metode til at gemme og hente data effektivt.

Software opbygges ofte i moduler med klare grænseflader til hinanden.

Datastrukturer

Datastruktur = metode til at gemme og hente data effektivt.

Software opbygges ofte i moduler med klare grænseflader til hinanden.

Datastrukturer kan beskrives på to niveauer:

Interface/grænseflade: Hvordan kan man tilgå data fra andre dele af et program (hvilke funktioner på data tilbydes)?

Implementation: Hvordan skal disse funktioner implementeres (hvordan skrives kode, som udfører dem)?

Datastrukturer

Datastruktur = metode til at gemme og hente data effektivt.

Software opbygges ofte i moduler med klare grænseflader til hinanden.

Datastrukturer kan beskrives på to niveauer:

Interface/grænseflade: Hvordan kan man tilgå data fra andre dele af et program (hvilke funktioner på data tilbydes)?

Implementation: Hvordan skal disse funktioner implementeres (hvordan skrives kode, som udfører dem)?

Dagens emne: to udbredte grundtyper (blandt mange flere) af interfaces:

▶ **Sekventiel tilgang**

▶ **Random access**

Interface for sekventiel tilgang

Interface for sekventiel tilgang

Funktioner til læsning:

```
readNext(), isEndOfFile(), open(), close()
```


Interface for sekventiel tilgang

Funktioner til læsning:

```
readNext(), isEndOfFile(), open(), close()
```

Typisk anvendelse:

```
file = open(filename)
```

```
Så længe file.isEndOfFile() svarer False:
```

```
    x = file.readNext()
```

```
    (...gør noget med x...)
```

```
file.close()
```

7	2	3	4	9	4	7	2	8	1	6
---	---	---	---	---	---	---	---	---	---	---

Interface for sekventiel tilgang

Funktioner til læsning:

```
readNext(), isEndOfFile(), open(), close()
```

Typisk anvendelse:

```
file = open(filename)
```

```
Så længe file.isEndOfFile() svarer False:
```

```
    x = file.readNext()
```

```
    (...gør noget med x...)
```

```
file.close()
```

7	2	3	4	9	4	7	2	8	1	6
---	---	---	---	---	---	---	---	---	---	---



Interface for sekventiel tilgang

Funktioner til læsning:

```
readNext(), isEndOfFile(), open(), close()
```

Typisk anvendelse:

```
file = open(filename)
```

```
Så længe file.isEndOfFile() svarer False:
```

```
    x = file.readNext()
```

```
    (...gør noget med x...)
```

```
file.close()
```

7	2	3	4	9	4	7	2	8	1	6
---	---	---	---	---	---	---	---	---	---	---



Interface for sekventiel tilgang

Funktioner til læsning:

```
readNext(), isEndOfFile(), open(), close()
```

Typisk anvendelse:

```
file = open(filename)
```

```
Så længe file.isEndOfFile() svarer False:
```

```
    x = file.readNext()
```

```
    (...gør noget med x...)
```

```
file.close()
```

7	2	3	4	9	4	7	2	8	1	6
---	---	---	---	---	---	---	---	---	---	---

↑ ...

Interface for sekventiel tilgang

Funktioner til læsning:

```
readNext(), isEndOfFile(), open(), close()
```

Typisk anvendelse:

```
file = open(filename)
```

```
Så længe file.isEndOfFile() svarer False:
```

```
    x = file.readNext()
```

```
    (...gør noget med x...)
```

```
file.close()
```

7	2	3	4	9	4	7	2	8	1	6
---	---	---	---	---	---	---	---	---	---	---

... ↑

Interface for sekventiel tilgang

Funktioner til læsning:

```
readNext(), isEndOfFile(), open(), close()
```

Typisk anvendelse:

```
file = open(filename)
```

```
Så længe file.isEndOfFile() svarer False:
```

```
    x = file.readNext()
```

```
    (...gør noget med x...)
```

```
file.close()
```

7	2	3	4	9	4	7	2	8	1	6
---	---	---	---	---	---	---	---	---	---	---



Interface for sekventiel tilgang

Funktioner til skrivning:

```
writeNext(), open(), close()
```

Interface for sekventiel tilgang

Funktioner til skrivning:

```
writeNext(), open(), close()
```



Interface for sekventiel tilgang

Funktioner til skrivning:

`writeNext()`, `open()`, `close()`

4



Interface for sekventiel tilgang

Funktioner til skrivning:

`writeNext()`, `open()`, `close()`

4	6
---	---



Interface for sekventiel tilgang

Funktioner til skrivning:

`writeNext()`, `open()`, `close()`

4	6	1
---	---	---



Interface for sekventiel tilgang

Funktioner til skrivning:

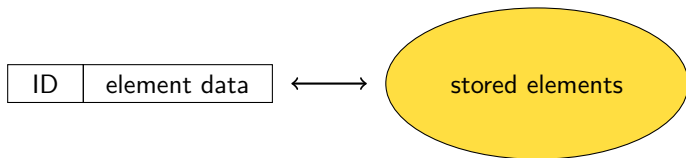
`writeNext()`, `open()`, `close()`

4	6	1	7
---	---	---	---



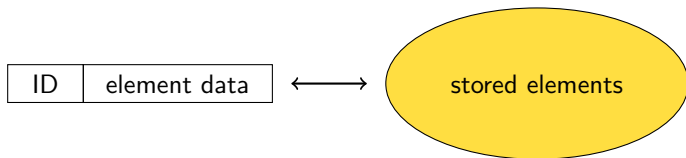
Interface for random access

Tilgang via ID (også kaldet key, nøgle) for elementer.



Interface for random access

Tilgang via ID (også kaldet key, nøgle) for elementer.



Funktioner:

`findElm(ID)` (returnerer element data)

`insertElm(ID,data)`

`deleteElm(ID)`

`open(), close()`

Interface for random access

```
findElm(ID)  
insertElm(ID,data)  
deleteElm(ID)  
open(), close()
```

Eksempler:

Interface for random access

```
findElm(ID)
insertElm(ID,data)
deleteElm(ID)
open(), close()
```

Eksempler:

- ▶ Databaser (ID = CPR, f.eks.)

Interface for random access

```
findElm(ID)
insertElm(ID,data)
deleteElm(ID)
open(), close()
```

Eksempler:

- ▶ Databaser (ID = CPR, f.eks.)
- ▶ Dictionaries i Python

```
data = myDict[ID]
myDict[ID] = data
myDict.pop(ID) (eller del myDict[ID])
```

Interface for random access

```
findElm(ID)
insertElm(ID,data)
deleteElm(ID)
open(), close()
```

Eksempler:

- ▶ Databaser (ID = CPR, f.eks.)
- ▶ Dictionaries i Python

```
data = myDict[ID]
myDict[ID] = data
myDict.pop(ID) (eller del myDict[ID])
```
- ▶ Lister i Python kan ses som et specialtilfælde, hvor ID = index

```
data = myList[7]
myList[7] = data
myList[7] = None
```

Dagens spørgsmål

1. Hvad kan det simple interface Sekventiel tilgang bruges til?

Dagens spørgsmål

1. Hvad kan det simple interface `Sekventiel tilgang` bruges til?
2. *Hvordan implementeres det mere avancerede interface `Random access`?*

Dagens spørgsmål

1. Hvad kan det simple interface `Sekventiel tilgang` bruges til?
2. *Hvordan implementeres det mere avancerede interface `Random access`?*

Bemærk: Alle datakilder kan tilgås med sekventiel tilgang. (Og for nogle kan *kun* sekventiel tilgang laves effektivt.)

- ▶ Harddisk
- ▶ SSD (Solid State Disk)
- ▶ CD
- ▶ Bånd
- ▶ Streaming over net
- ▶ Data genereret on-the-fly af et andet program
- ▶ Data i en liste (Python) eller i et array (Java).

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element
- ▶ Find sum og antal af elementer (og dermed gennemsnit)

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element
- ▶ Find sum og antal af elementer (og dermed gennemsnit)
- ▶ Fjern mindste element

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element
- ▶ Find sum og antal af elementer (og dermed gennemsnit)
- ▶ Fjern mindste element
- ▶ Selectionsort (*detaljer om lidt*)

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element
- ▶ Find sum og antal af elementer (og dermed gennemsnit)
- ▶ Fjern mindste element
- ▶ Selectionsort (*detaljer om lidt*)
- ▶ Merge: slå to sortererede lister sammen til én sorteret liste (*detaljer om lidt*)

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element
- ▶ Find sum og antal af elementer (og dermed gennemsnit)
- ▶ Fjern mindste element
- ▶ Selectionsort (*detaljer om lidt*)
- ▶ Merge: slå to sortererede lister sammen til én sorteret liste (*detaljer om lidt*)
- ▶ Mergesort (*detaljer om lidt*)

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element
- ▶ Find sum og antal af elementer (og dermed gennemsnit)
- ▶ Fjern mindste element
- ▶ Selectionsort (*detaljer om lidt*)
- ▶ Merge: slå to sortererede lister sammen til én sorteret liste (*detaljer om lidt*)
- ▶ Mergesort (*detaljer om lidt*)
- ▶ Generaliseret merge (*detaljer om lidt*)

Spørgsmål 1: Hvad kan laves med sekventiel tilgang?

F.eks.:

- ▶ Lineær søgning (jvf. Lenes slides)
- ▶ Find mindste element
- ▶ Find sum og antal af elementer (og dermed gennemsnit)
- ▶ Fjern mindste element
- ▶ Selectionsort (*detaljer om lidt*)
- ▶ Merge: slå to sorterede lister sammen til én sorteret liste (*detaljer om lidt*)
- ▶ Mergesort (*detaljer om lidt*)
- ▶ Generaliseret merge (*detaljer om lidt*)

Morale: Man kan lave forbavsende mange algoritmer baseret på det simple interface Sekventiel tilgang.

Selectionsort

Først find mindste element via et scan af input:

7	3	5	8	5	9	2	4	6	4	2	7
---	---	---	---	---	---	---	---	---	---	---	---

I dette tilfælde er mindste element 2.

Selectionsort

Først find mindste element via et scan af input:

7	3	5	8	5	9	2	4	6	4	2	7
---	---	---	---	---	---	---	---	---	---	---	---

I dette tilfælde er mindste element 2. Udtag nu dette fra listen via et scan mere, hvor det første 2 skrives til output, mens alle andre elementer kopieres til en ny liste:

2	7	3	5	8	5	9	4	6	4	2	7
---	---	---	---	---	---	---	---	---	---	---	---

Selectionsort

Først find mindste element via et scan af input:

7	3	5	8	5	9	2	4	6	4	2	7
---	---	---	---	---	---	---	---	---	---	---	---

I dette tilfælde er mindste element 2. Udtag nu dette fra listen via et scan mere, hvor det første 2 skrives til output, mens alle andre elementer kopieres til en ny liste:

2	7	3	5	8	5	9	4	6	4	2	7
---	---	---	---	---	---	---	---	---	---	---	---

Gentag nu ovenstående. Efter seks runder vil vi f.eks. have følgende situation:

2	2	3	4	4	5	7	8	5	9	6	7
---	---	---	---	---	---	---	---	---	---	---	---

Selectionsort

Først find mindste element via et scan af input:

7	3	5	8	5	9	2	4	6	4	2	7
---	---	---	---	---	---	---	---	---	---	---	---

I dette tilfælde er mindste element 2. Udtag nu dette fra listen via et scan mere, hvor det første 2 skrives til output, mens alle andre elementer kopieres til en ny liste:

2

7	3	5	8	5	9	4	6	4	2	7
---	---	---	---	---	---	---	---	---	---	---

Gentag nu ovenstående. Efter seks runder vil vi f.eks. have følgende situation:

2	2	3	4	4	5
---	---	---	---	---	---

7	8	5	9	6	7
---	---	---	---	---	---

Efter n runder vil input være blevet sorteret:

2	2	3	4	4	5	5	6	7	7	8	9
---	---	---	---	---	---	---	---	---	---	---	---

Selectionsort

Korrekthed af Selectionsort er oplagt (vi udtager først det mindste element, så det næstmindste, så det tredjemindste, . . .).

Selectionsort

Korrekthed af Selectionsort er oplagt (vi udtager først det mindste element, så det næstmindste, så det trediemindste, . . .).

Køretid for Selectionsort?

Selectionsort

Korrekthed af Selectionsort er oplagt (vi udtager først det mindste element, så det næstmindste, så det tredjemindste, ...).

Køretid for Selectionsort? Den er proportional med:

$$n + (n - 1) + (n - 2) + \dots + 3 + 2 + 1$$

Selectionsort

Korrekthed af Selectionsort er oplagt (vi udtager først det mindste element, så det næstmindste, så det trediemindste, ...).

Køretid for Selectionsort? Den er proportional med:

$$n + (n - 1) + (n - 2) + \dots + 3 + 2 + 1$$

Kald ovenstående sum for S . Så gælder

$$2S = \begin{array}{cccccccc} n & + & (n-1) & + & (n-2) & + & \dots & + & 3 & + & 2 & + & 1 \\ + & 1 & + & 2 & + & 3 & + & \dots & + & (n-2) & + & (n-1) & + & n \end{array}$$

Selectionsort

Korrekthed af Selectionsort er oplagt (vi udtager først det mindste element, så det næstmindste, så det tredjemindste, ...).

Køretid for Selectionsort? Den er proportional med:

$$n + (n - 1) + (n - 2) + \dots + 3 + 2 + 1$$

Kald ovenstående sum for S . Så gælder

$$\begin{aligned} 2S = & \quad n + (n - 1) + (n - 2) + \dots + 3 + 2 + 1 \\ & + 1 + 2 + 3 + \dots + (n - 2) + (n - 1) + n \end{aligned}$$

Dvs. $2S = n(n + 1)$, hvoraf vi ser $S = n(n + 1)/2 = n^2/2 + n/2$.

Selectionsort

Korrekthed af Selectionsort er oplagt (vi udtager først det mindste element, så det næstmindste, så det trediemindste, ...).

Køretid for Selectionsort? Den er proportional med:

$$n + (n - 1) + (n - 2) + \dots + 3 + 2 + 1$$

Kald ovenstående sum for S . Så gælder

$$2S = \begin{array}{cccccccc} n & + & (n-1) & + & (n-2) & + & \dots & + & 3 & + & 2 & + & 1 \\ + & 1 & + & 2 & + & 3 & + & \dots & + & (n-2) & + & (n-1) & + & n \end{array}$$

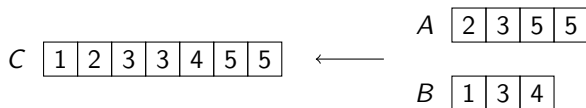
Dvs. $2S = n(n+1)$, hvoraf vi ser $S = n(n+1)/2 = n^2/2 + n/2$.

Da n^2 dominerer (vokser hurtigere end) n , er køretiden

$$O(n^2)$$

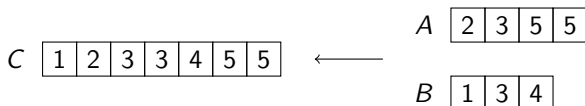
Merge

Mål: Slå to sorterede lister A og B sammen til én sorteret liste C .



Merge

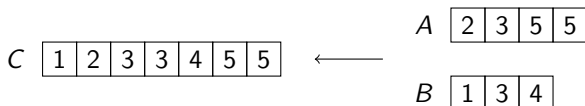
Mål: Slå to sorterede lister A og B sammen til én sorteret liste C .



Et **merge-skridt**: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Merge

Mål: Slå to sorterede lister A og B sammen til én sorteret liste C .



Et **merge-skridt**: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Start med tom C

Så længe både A og B er ikke-tomme:

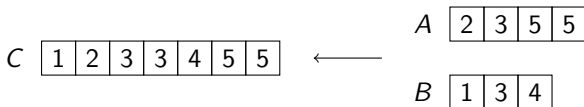
 Udfør et merge-skridt

Hvis enten A eller B er ikke-tom:

 Flyt resten af dens elementer over sidst i C

Merge

Mål: Slå to sorterede lister A og B sammen til én sorteret liste C .



Et **merge-skridt**: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Start med tom C

Så længe både A og B er ikke-tomme:

 Udfør et merge-skridt

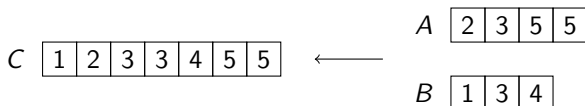
Hvis enten A eller B er ikke-tom:

 Flyt resten af dens elementer over sidst i C

Er det en korrekt metode (dvs. er C ved afslutning sorteret)?

Merge

Mål: Slå to sorterede lister A og B sammen til én sorteret liste C .



Et **merge-skridt**: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Start med tom C

Så længe både A og B er ikke-tomme:

 Udfør et merge-skridt

Hvis enten A eller B er ikke-tom:

 Flyt resten af dens elementer over sidst i C

Er det en korrekt metode (dvs. er C ved afslutning sorteret)?

Korrekthed kan ses ved hjælp af en **invariant**.

Korrekthed og køretid for merge

Mål: Slå to (stigende) sorterede lister A og B sammen til én sorteret liste C .

Et merge-skridt: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Start med tom C

Så længe både A og B er ikke-tomme:

 Udfør et merge-skridt

Hvis enten A eller B er ikke-tom:

 Flyt resten af dens elementer over sidst i C

Korrekthed og køretid for merge

Mål: Slå to (stigende) sorterede lister A og B sammen til én sorteret liste C .

Et merge-skridt: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Start med tom C

Så længe både A og B er ikke-tomme:

 Udfør et merge-skridt

Hvis enten A eller B er ikke-tom:

 Flyt resten af dens elementer over sidst i C

Korrekthed (at C ved afslutning er sorteret) kan ses ved flg. **invariant**:

Nuværende version af A , B , C er sorterede, og ingen elementer i A , B er mindre end noget element i C .

Korrekthed og køretid for merge

Mål: Slå to (stigende) sorterede lister A og B sammen til én sorteret liste C .

Et merge-skridt: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Start med tom C

Så længe både A og B er ikke-tomme:

 Udfør et merge-skridt

Hvis enten A eller B er ikke-tom:

 Flyt resten af dens elementer over sidst i C

Korrekthed (at C ved afslutning er sorteret) kan ses ved flg. **invariant**:

Nuværende version af A , B , C er sorterede, og ingen elementer i A , B er mindre end noget element i C .

Tid:

Korrekthed og køretid for merge

Mål: Slå to (stigende) sorterede lister A og B sammen til én sorteret liste C .

Et merge-skridt: Sammenlign (nuværende) forreste element i A og B , og flyt det mindste af disse over sidst i C .

Start med tom C

Så længe både A og B er ikke-tomme:

 Udfør et merge-skridt

Hvis enten A eller B er ikke-tom:

 Flyt resten af dens elementer over sidst i C

Korrekthed (at C ved afslutning er sorteret) kan ses ved flg. **invariant**:

Nuværende version af A , B , C er sorterede, og ingen elementer i A , B er mindre end noget element i C .

Tid: $O(|A| + |B|)$, da hvert element røres én gang.

Mergesort

Runde 1: Merge lister af længde 1 (disse er automatisk sorterede) parvis sammen til sorterede lister af længde 2 ($= 2^1$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 2: Merge de sorterede lister af længde 2 til sorterede lister af længde 4 ($= 2^2$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 3...: Længde 4 til længde 8 ($= 2^3$),...

Efter i runder: længde $= 2^i$.

Indtil længde n nås, dvs. alle elementer er i én sorteret liste. Dvs. vi har sorteret elementerne.

Mergesort

Runde 1: Merge lister af længde 1 (disse er automatisk sorterede) parvis sammen til sorterede lister af længde 2 ($= 2^1$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 2: Merge de sorterede lister af længde 2 til sorterede lister af længde 4 ($= 2^2$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 3...: Længde 4 til længde 8 ($= 2^3$),...

Efter i runder: længde $= 2^i$.

Indtil længde n nås, dvs. alle elementer er i én sorteret liste. Dvs. vi har sorteret elementerne.

Der er lavet $\log_2 n$ runder, når længde n nås [da $2^i = n \Leftrightarrow i = \log_2 n$].

Mergesort

Runde 1: Merge lister af længde 1 (disse er automatisk sorterede) parvis sammen til sorterede lister af længde 2 ($= 2^1$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 2: Merge de sorterede lister af længde 2 til sorterede lister af længde 4 ($= 2^2$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 3...: Længde 4 til længde 8 ($= 2^3$),...

Efter i runder: længde $= 2^i$.

Indtil længde n nås, dvs. alle elementer er i én sorteret liste. Dvs. vi har sorteret elementerne.

Der er lavet $\log_2 n$ runder, når længde n nås [da $2^i = n \Leftrightarrow i = \log_2 n$].

Tid i alt:

Mergesort

Runde 1: Merge lister af længde 1 (disse er automatisk sorterede) parvis sammen til sorterede lister af længde 2 ($= 2^1$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 2: Merge de sorterede lister af længde 2 til sorterede lister af længde 4 ($= 2^2$). Arbejdet svarer til at røre alle elementer én gang, dvs. er $O(n)$.

Runde 3...: Længde 4 til længde 8 ($= 2^3$),...

Efter i runder: længde $= 2^i$.

Indtil længde n nås, dvs. alle elementer er i én sorteret liste. Dvs. vi har sorteret elementerne.

Der er lavet $\log_2 n$ runder, når længde n nås [da $2^i = n \Leftrightarrow i = \log_2 n$].

Tid i alt: $O(n \log_2 n)$ [da vi $\log_2 n$ gange laver $O(n)$ arbejde].

Generaliseret merge

Under merge af to sorterede lister A og B vil elementer med samme værdi “møde” hinanden under merge-skridtene (dvs. de vil på et tidspunkt stå forrest samtidigt i A og B). Dette kan udnyttes:

Generaliseret merge

Under merge af to sorterede lister A og B vil elementer med samme værdi “møde” hinanden under merge-skridtene (dvs. de vil på et tidspunkt stå forrest samtidigt i A og B). Dette kan udnyttes:

Eksempel 1: A er datafil med (ID,data)-elementer, sorteret efter ID. B er liste af (ID,dataopdatering)-elementer, sorteret efter ID.

At gennemløbe A og B via merge-skridt vil tillade os at opdatere data i A . NB: de to lister behøver ikke være lige lange.

Generaliseret merge

Under merge af to sorterede lister A og B vil elementer med samme værdi “møde” hinanden under merge-skridtene (dvs. de vil på et tidspunkt stå forrest samtidigt i A og B). Dette kan udnyttes:

Eksempel 1: A er datafil med (ID,data)-elementer, sorteret efter ID. B er liste af (ID,dataopdatering)-elementer, sorteret efter ID.

At gennemløbe A og B via merge-skridt vil tillade os at opdatere data i A . NB: de to lister behøver ikke være lige lange.

Eksempel 2: To sorterede lister A og B repræsenterer mængder (og er derfor hver især uden dubletter).

At gennemløbe A og B via merge-skridt vil tillade os at generere f.eks. $A \cap B$: Hvis et merge-step ser to ens elementer som forreste i A og B , flyttes det ene over sidst i C , og det andet smides væk. Hvis et merge-step ser to forskellige, smides det mindste væk.

Spørgsmål 2: Hvordan laves interface Random access?

Én metode: Hashing.

Spørgsmål 2: Hvordan laves interface Random access?

Én metode: Hashing.



Spørgsmål 2: Hvordan laves interface Random access?

Én metode: Hashing.



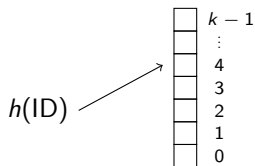
Hash \sim hacher [fransk] \sim hakke i stykker.

Hashing

Ide: Tildel hver ID et index i i et array (Python: en liste) A , hvor det tilhørende element gemmes.

Hash-funktion h :

$$h(\text{ID}) = \text{index } i \text{ i } A$$



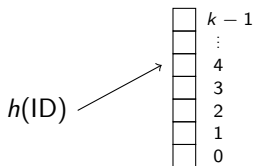
Hashing

Ide: Tildel hver ID et index i i et array (Python: en liste) A , hvor det tilhørende element gemmes.

Hash-funktion h :

$$h(\text{ID}) = \text{index i } A$$

Eksempel (antag ID'er er heltal):



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

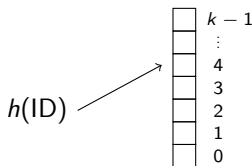
hvor $k = |A|$ og mod (modulus) betegner rest ved heltalsdivision.

Hashing

Ide: Tildel hver ID et index i i et array (Python: en liste) A , hvor det tilhørende element gemmes.

Hash-funktion h :

$$h(\text{ID}) = \text{index } i \text{ i } A$$



Eksempel (antag ID'er er heltal):

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

hvor $k = |A|$ og mod (modulus) betegner rest ved heltalsdivision.

Med $k = 41$:	$h(46) = 5$	(da $1 \cdot 41 + 5 = 46$)
	$h(12) = 12$	(da $0 \cdot 41 + 12 = 12$)
	$h(100) = 18$	(da $2 \cdot 41 + 18 = 100$)
	$h(479869) = 5$	(da $11704 \cdot 41 + 5 = 479869$)

Bemærk at $h(x) \in \{0, 1, 2, \dots, k-1\}$, så det er altid et lovligt index.

Hashing

Når nu vi antager, at ID'er er heltal, hvorfor så ikke bare bruge

$$h(x) = x?$$

Hashing

Når nu vi antager, at ID'er er heltal, hvorfor så ikke bare bruge

$$h(x) = x?$$

Eksempel: gem 5 CPR-numre.

Hashing

Når nu vi antager, at ID'er er heltal, hvorfor så ikke bare bruge

$$h(x) = x?$$

Eksempel: gem 5 CPR-numre.

CPR-numre: $180781-2345 \in \{0, 1, 2, \dots, 10^{10} - 1\}$

Så med ovenstående $h(x)$ skal størrelsen k af A skal være 10^{10} for at gemme 5 tal. Spild af plads (10^{10} bytes = 10 Gb).

Hashing

Når nu vi antager, at ID'er er heltal, hvorfor så ikke bare bruge

$$h(x) = x?$$

Eksempel: gem 5 CPR-numre.

CPR-numre: $180781-2345 \in \{0, 1, 2, \dots, 10^{10} - 1\}$

Så med ovenstående $h(x)$ skal størrelsen k af A skal være 10^{10} for at gemme 5 tal. Spild af plads (10^{10} bytes = 10 Gb).

Ofte er nøgler heltal (32 eller 64 bits), dvs. har $2^{32} \approx 0.43 \cdot 10^{10}$ eller $2^{64} \approx 0.18 \cdot 10^{20}$ muligheder. Dvs. samme situation eller værre.

Hashing

Når nu vi antager, at ID'er er heltal, hvorfor så ikke bare bruge

$$h(x) = x?$$

Eksempel: gem 5 CPR-numre.

CPR-numre: $180781-2345 \in \{0, 1, 2, \dots, 10^{10} - 1\}$

Så med ovenstående $h(x)$ skal størrelsen k af A skal være 10^{10} for at gemme 5 tal. Spild af plads (10^{10} bytes = 10 Gb).

Ofte er nøgler heltal (32 eller 64 bits), dvs. har $2^{32} \approx 0.43 \cdot 10^{10}$ eller $2^{64} \approx 0.18 \cdot 10^{20}$ muligheder. Dvs. samme situation eller værre.

Generelt ønsker vi, at længden k af A cirka er lig antal gemte elementer n (så hashtabellen selv ikke fylder væsentlig mere end de gemte elementer).

Dette opnås ved at vælge k i nærheden af n i udtrykket $h(x) = x \bmod k$. [Om nødvendigt må hashtabellen genopbygges med større k , hvis n vokser for meget pga. indsættelser.]

Kollisioner

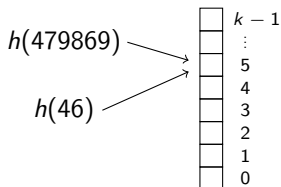
$h(x) = x \bmod 41$:

$h(46) = 5$	(da $1 \cdot 41 + 5 = 46$)
$h(12) = 12$	(da $0 \cdot 41 + 12 = 12$)
$h(100) = 18$	(da $2 \cdot 41 + 18 = 100$)
$h(479869) = 5$	(da $11704 \cdot 41 + 5 = 479869$)

Kollisioner

$$h(x) = x \bmod 41:$$

$$\begin{aligned}h(46) &= 5 && \text{(da } 1 \cdot 41 + 5 = 46\text{)} \\h(12) &= 12 && \text{(da } 0 \cdot 41 + 12 = 12\text{)} \\h(100) &= 18 && \text{(da } 2 \cdot 41 + 18 = 100\text{)} \\h(479869) &= 5 && \text{(da } 11704 \cdot 41 + 5 = 479869\text{)}\end{aligned}$$



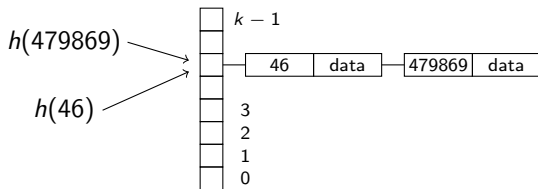
Problem: flere elementer vil bo på samme plads i A .

Én løsning: Chaining

Gem alle elementer på samme plads i en liste.

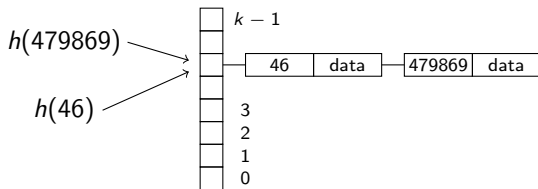
Én løsning: Chaining

Gem alle elementer på samme plads i en liste.



Én løsning: Chaining

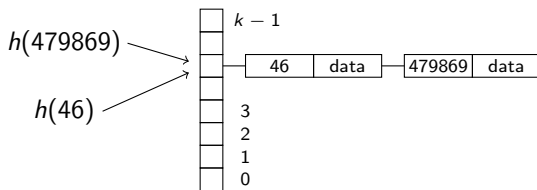
Gem alle elementer på samme plads i en liste.



De listen skal gennemløbes sekventielt (under søgning og under sletning), stiger tiden fra $O(1)$ til $O(|\text{liste}|)$.

Én løsning: Chaining

Gem alle elementer på samme plads i en liste.



De listen skal gennemløbes sekventielt (under søgning og under sletning), stiger tiden fra $O(1)$ til $O(|\text{liste}|)$.

Vi vil derfor gerne have få kollisioner.

Kan kollisioner undgås?

Det afhænger af forholdet mellem hash-funktionen h og det konkrete datasæt, som er gemt.

Kan kollisioner undgås?

Det afhænger af forholdet mellem hash-funktionen h og det konkrete datasæt, som er gemt.

Nogle observationer:

- ▶ I værste fald hash'er alle n elementer til samme plads. Tid: $O(n)$.

Kan kollisioner undgås?

Det afhænger af forholdet mellem hash-funktionen h og det konkrete datasæt, som er gemt.

Nogle observationer:

- ▶ I værste fald hash'er alle n elementer til samme plads. Tid: $O(n)$.
- ▶ Hvis n (antal elementer gemt) er større end k (størrelsen af A), er der mindst én kollision

Kan kollisioner undgås?

Det afhænger af forholdet mellem hash-funktionen h og det konkrete datasæt, som er gemt.

Nogle observationer:

- ▶ I værste fald hash'er alle n elementer til samme plads. Tid: $O(n)$.
- ▶ Hvis n (antal elementer gemt) er større end k (størrelsen af A), er der mindst én kollision (duehulsprincippet).

Kan kollisioner undgås?

Det afhænger af forholdet mellem hash-funktionen h og det konkrete datasæt, som er gemt.

Nogle observationer:

- ▶ I værste fald hash'er alle n elementer til samme plads. Tid: $O(n)$.
- ▶ Hvis n (antal elementer gemt) er større end k (størrelsen af A), er der mindst én kollision (duehulsprincippet).

Hvad hvis vi antager at h "indsætter tallene på tilfældige pladser"?

NB: Dette er ikke en brugbar hash-funktion, da vi ikke kan finde elementerne igen (vi kan ikke huske, hvor de er lagt). Men det giver en jævn spredning i hashtabellen for *alle* datasæt, og kan derfor ses som en opførelse, som vi ønsker for en hash-funktion.

Kan kollisioner undgås?

Det afhænger af forholdet mellem hash-funktionen h og det konkrete datasæt, som er gemt.

Nogle observationer:

- ▶ I værste fald hash'er alle n elementer til samme plads. Tid: $O(n)$.
- ▶ Hvis n (antal elementer gemt) er større end k (størrelsen af A), er der mindst én kollision (duehulsprincippet).

Hvad hvis vi antager at h "indsætter tallene på tilfældige pladser"?

NB: Dette er ikke en brugbar hash-funktion, da vi ikke kan finde elementerne igen (vi kan ikke huske, hvor de er lagt). Men det giver en jævn spredning i hashtabellen for *alle* datasæt, og kan derfor ses som en opførelse, som vi ønsker for en hash-funktion.

Vi studerer derfor nu denne situation, og ønsker at sige noget om, *hvor hurtigt* kollisioner opstår, når vi indsætter elementer tilfældigt.

Fødselsdage og hashing

Fødselsdage og hashing

Situation: n tilfældige personer går ind i et rum.

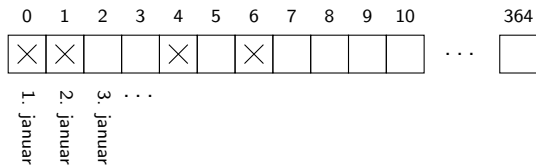
Fokus: Er der nogen, som har fødselsdag samme dato i året?

Fødselsdage og hashing

Situation: n tilfældige personer går ind i et rum.

Fokus: Er der nogen, som har fødselsdag samme dato i året?

Bemærk: at samle personer med tilfældige fødselsdage svarer til at indsætte tilfældigt i en hashtabel af størrelse 365 (= antal dage i et år).



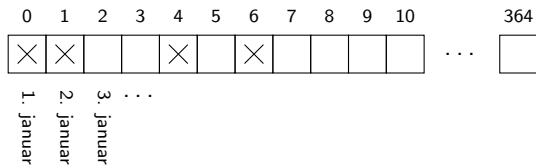
Ovenfor er fire personer samlet (hvoraf ingen har fødselsdag samme dag).

Fødselsdage og hashing

Situation: n tilfældige personer går ind i et rum.

Fokus: Er der nogen, som har fødselsdag samme dato i året?

Bemærk: at samle personer med tilfældige fødselsdage svarer til at indsætte tilfældigt i en hashtabel af størrelse 365 (= antal dage i et år).



Ovenfor er fire personer samlet (hvoraf ingen har fødselsdag samme dag).

Dette er altså et eksempel på den (hashing) situation vi ønsker at studere.

Fødselsdage og hashing

Situation: n tilfældige personer går ind i et rum.

Fokus: Er der nogen, som har fødselsdag samme dato i året?

Fødselsdage og hashing

Situation: n tilfældige personer går ind i et rum.

Fokus: Er der nogen, som har fødselsdag samme dato i året?

Sandsynligheden for dette må stige, når n vokser:

n	Sandsynlighed for (mindst) to med samme fødselsdag
0	0
1	0
2	$1/365$
\vdots	\vdots
?	$1/2$
\vdots	\vdots
366	1

Spørgsmål: For hvilket n bliver sandsynligheden for, at **nogen har med samme fødselsdag, større end $1/2$?**

Fødselsdage og hashing

Spørgsmål: For hvilket n bliver sandsynligheden for, at nogen har samme fødselsdag, større end $1/2$?

Bemærk:

$$1 - P(\text{nogen har samme fødseldag}) = P(\text{ingen har samme fødseldag})$$

Fødselsdage og hashing

Spørgsmål: For hvilket n bliver sandsynligheden for, at **nogen har samme fødselsdag**, større end $1/2$?

Bemærk:

$$1 - P(\text{nogen har samme fødseldag}) = P(\text{ingen har samme fødseldag})$$

Så for at svare på spørgsmålet kigger vi i stedet på følgende sandsynlighed

$$P(\text{ingen med samme fødseldag blandt de } n \text{ første personer}),$$

og spørger, hvornår den er *mindre* end $1 - 1/2 = 1/2$.

Fødselsdage og hashing

Ingen med samme fødseldag blandt de n første personer



- 1) Ingen med samme fødselsdag blandt de $n - 1$ første personer
og
- 2) Den n 'te persons fødselsdag er ikke lig nogen af disse.

Fødselsdage og hashing

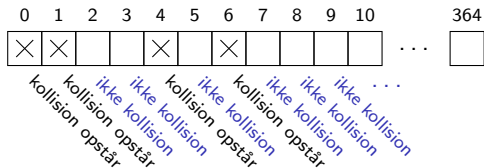
Ingen med samme fødseldag blandt de n første personer



- 1) Ingen med samme fødselsdag blandt de $n - 1$ første personer
og
- 2) Den n 'te persons fødselsdag er ikke lig nogen af disse.

Hvis 1) gælder, er præcis $n - 1$ forskellige datoer optaget, når den n 'te person går ind i rummet. Hvis personerne antages at have tilfældige fødselsdage, er sandsynligheden for 2) derfor $(365 - (n - 1))/365$.

Situationen for $n = 5$ (dvs. 4 personer allerede til stede, uden kollision):



Fødselsdage og hashing

Hvis vi kalder P (ingen med samme fødseldag blandt de n første personer) for s_n , kan vi se af ovenstående at

$$s_n = s_{n-1} \cdot \frac{365 - (n - 1)}{365}$$

Da s_1 naturligvis er 1 (med kun én person i rummet er der helt sikkert ingen med samme fødselsdag), ser vi at:

$$s_1 = 1$$

$$s_2 = s_1 \cdot \frac{364}{365} = 1 \cdot \frac{364}{365} = 0.9972 \dots$$

$$s_3 = s_2 \cdot \frac{363}{365} = 1 \cdot \frac{364}{365} \cdot \frac{363}{365} = 0.9917 \dots$$

$$s_4 = s_3 \cdot \frac{362}{365} = 1 \cdot \frac{364}{365} \cdot \frac{363}{365} \cdot \frac{362}{365} = 0.9836 \dots$$

⋮

Fødselsdagsparadokset

Hvis man udregner disse udtryk (evt. via at skrive et program, der gør det), ses at

$$s_{22} = 0.5243 \dots$$

$$s_{23} = 0.4927 \dots$$

Fødselsdagsparadokset

Hvis man udregner disse udtryk (evt. via at skrive et program, der gør det), ses at

$$s_{22} = 0.5243 \dots$$

$$s_{23} = 0.4927 \dots$$

Vi var interesserede i, for hvilket n det var mere sandsynligt, at nogen havde samme fødselsdags, end at ingen havde.

Dvs. for hvilket n at $s_n \leq 1/2$.

Af beregningen ovenfor får vi vores svar: allerede når n bliver 23.

Fødselsdagsparadokset

Hvis man udregner disse udtryk (evt. via at skrive et program, der gør det), ses at

$$s_{22} = 0.5243 \dots$$

$$s_{23} = 0.4927 \dots$$

Vi var interesserede i, for hvilket n det var mere sandsynligt, at nogen havde samme fødselsdags, end at ingen havde.

Dvs. for hvilket n at $s_n \leq 1/2$.

Af beregningen ovenfor får vi vores svar: allerede når n bliver 23.

Navnet “fødselsdagsparadokset” bruges, fordi dette er et meget mindre antal personer end de fleste intuitivt vil gætte på.

Fødselsdagsparadokset

Hvis man udregner disse udtryk (evt. via at skrive et program, der gør det), ses at

$$s_{22} = 0.5243 \dots$$

$$s_{23} = 0.4927 \dots$$

Vi var interesserede i, for hvilket n det var mere sandsynligt, at nogen havde samme fødselsdags, end at ingen havde.

Dvs. for hvilket n at $s_n \leq 1/2$.

Af beregningen ovenfor får vi vores svar: allerede når n bliver 23.

Navnet “fødselsdagsparadokset” bruges, fordi dette er et meget mindre antal personer end de fleste intuitivt vil gætte på.

Af samme slags beregninger fås f.eks. $s_{50} = 0.0296 \dots$. Dvs. at der er under 3% sandsynlighed for, at der i en gruppe af 50 personer ikke er to med samme fødselsdato.

Fødselsdagsparadokset og hashing

Recall: At samle personer med tilfældige fødselsdage svarer til at indsætte tilfældigt i en hashtabel af størrelse 365 (antal dage i et år).

Ovenstående metode kan derfor bruges til at beregne sandsynligheden for, at kollisioner opstår ved indsættelse af n elementer i en hashtabel, hvis det antages, at hashfunktionen tildeler tabel-indekser til elementer på en tilfældig måde. Man skal blot udskifte 365 med den aktuelle tabelstørrelse.

Fødselsdagsparadokset og hashing

Recall: At samle personer med tilfældige fødselsdage svarer til at indsætte tilfældigt i en hashtabel af størrelse 365 (antal dage i et år).

Ovenstående metode kan derfor bruges til at beregne sandsynligheden for, at kollisioner opstår ved indsættelse af n elementer i en hashtabel, hvis det antages, at hashfunktionen tildeler tabel-indeksers til elementer på en tilfældig måde. Man skal blot udskifte 365 med den aktuelle tabelstørrelse.

Morale: vi kan se af sådanne beregninger, at den første kollision opstår overraskende hurtigt, når der indsættes i hashtabeller (selv hvis vi finder en hashfunktion, der opfører sig, som vi ønsker os).

Det er derfor nødvendigt at have et system til at klare kollisioner (f.eks. via lister, jvf. side 21 fra tidligere) for at kunne bruge hashing.

Opsummering om hashing

Vi har set:

- ▶ Hovedideen i hashing: Udnyt at heltal kan bruges som indekser til hukommelsesceller.
- ▶ Vi har brug for en funktion til at mappe store heltal (ID'er) til små heltal for ikke at spilde plads (vi vil gerne kunne vælge $k \approx n$).
- ▶ Det er nødvendigt at have et system til at klare kollisioner.

Opsummering om hashing

Vi har set:

- ▶ Hovedideen i hashing: Udnyt at heltal kan bruges som indekser til hukommelsesceller.
- ▶ Vi har brug for en funktion til at mappe store heltal (ID'er) til små heltal for ikke at spille plads (vi vil gerne kunne vælge $k \approx n$).
- ▶ Det er nødvendigt at have et system til at klare kollisioner.

Vi har ikke set:

- ▶ Hvordan man i praksis vælger en hash-funktion, som må forventes for de fleste datasæt at lave få kollisioner.
- ▶ Hvordan man laver hashfunktioner på IDs, som ikke er heltal (f.eks. strenge eller generelle objekter).
- ▶ Andre systemer (end lister) til at klare kollisioner.
- ▶ Andre anvendelser af hash-funktioner (end til at implementere datastrukturer for random access API'et): fingerprinting, randomisering af data, kryptografiske anvendelser.

Opsummering om hashing

I praksis ses ofte en todelt hash-funktion:

- ▶ En separat del (`hash()` i Python, `hashCode()` i Java) som tager en vilkårlig ID og omformer den til et 64 bit heltal.
- ▶ En del indbygget i random access datastrukturen (dictionaries i Python, `HashMap` i Java) som reducerer til et heltal mindre end k .

Med en (evt. todelt) hashfunktion, der er valgt med omhu (som i Python og Java), fungerer hashing meget effektivt i praksis: de fleste kald til datastrukturen vil tage $O(1)$ tid.