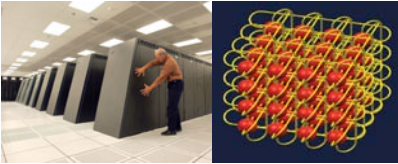
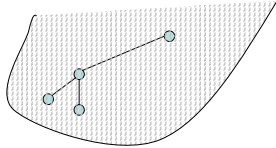


Algoritmer i distribuerte nettverk

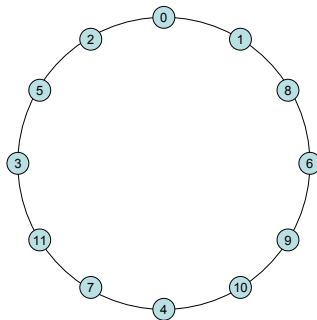
Parallele datamaskiner <small>Blue Gene/L, BGW, ASC Purple</small>	Distribuerte nettverk
<ul style="list-style-type: none"> • Topologien kjent. • Alle kan kommunisere med alle, "direkte". <p>Blue Gene/L (IBM) (Lawrence Livermore)</p> <p>131.072 prosessorer 367.000 GFlops (367 billioner flyttallsoperasjoner/sekund)</p> 	<ul style="list-style-type: none"> • Topologien ikke kjent – en prosessor kjenner bare naboene sine. • En melding kan bare sendes til en nabo.  <p>LAN, Internet, mobile ad-hoc nett, trådløse nett, etc.</p>

[Top 500 supercomputer sites](#)

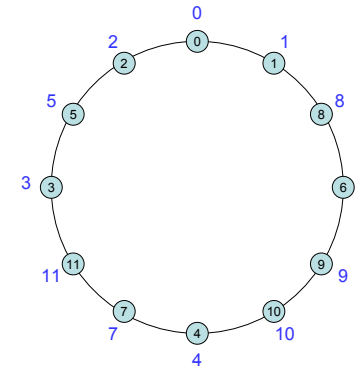
Algoritmer i distribuerte nettverk

- Topologien ikke kjent:
 - kjenner ikke strukturen,
 - kjenner ikke antall prosessorer.
- Ingen sentralisert kontroll.
- Som regel grafproblemer vi ser på:
 - nettverket er inndata for algoritmen.
- Trenger rutiner for å legge struktur på nettet vårt.
- Algoritmene går i steg (synkron modell):
 - alle prosessorene gjennomfører en blokk med kode, før de går videre.
- Tidskompleksitet (antall steg).
- Kommunikasjonskompleksitet (antall meldinger).

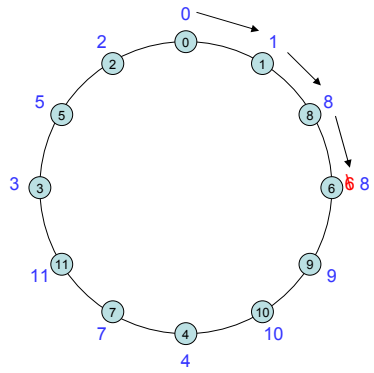
Ledervalg



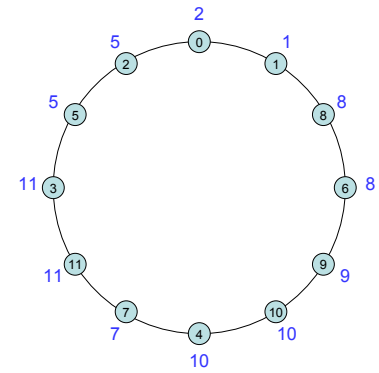
Ledervalg



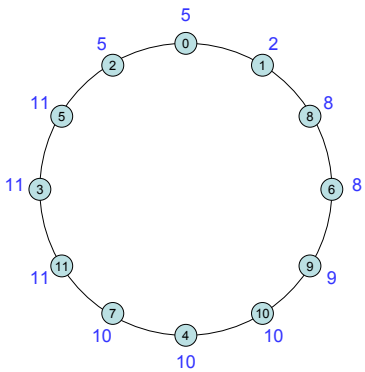
Ledervalg



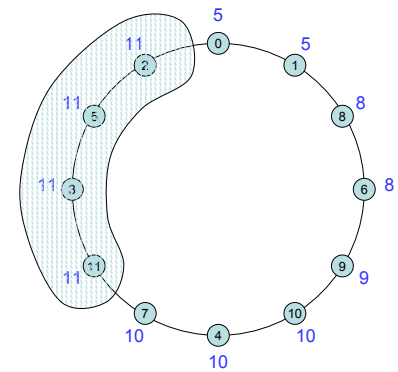
Ledervalg



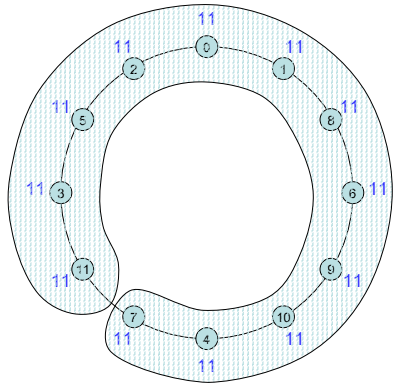
Ledervalg



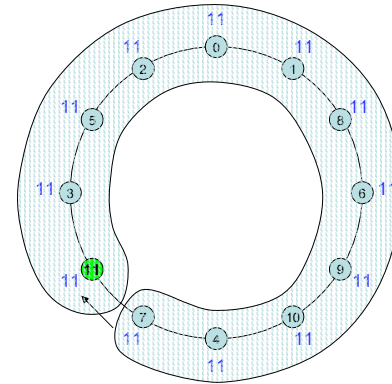
Ledervalg



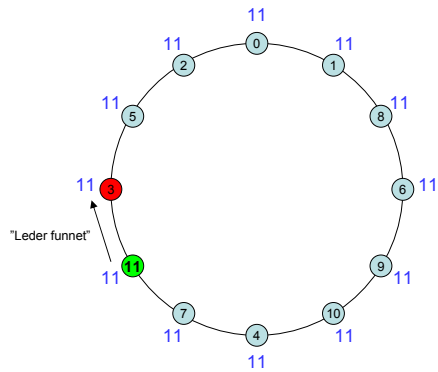
Ledervalg



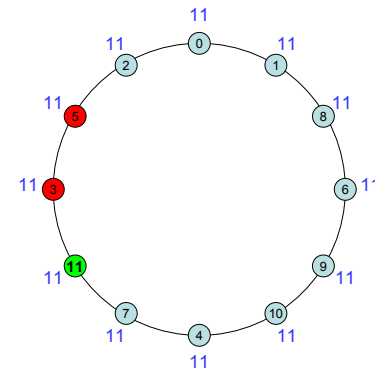
Ledervalg



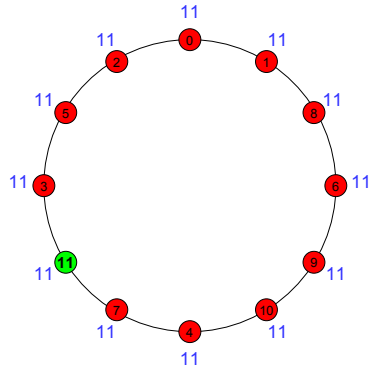
Ledervalg



Ledervalg



Ledervalg



Ledervalg

```

proc RingLeaderElection
  forall v do
    send v.UID to clockwise neighbour
  endforall

  while not < all v have received "leder found" message > do
    if u.UID received from counterclockwise neighbour of v then
      if v.UID < u.UID then
        send u.UID to clockwise neighbour
      else
        if v.UID = u.UID then
          v.leaderStatus = "leader"
          send "leder found" to clockwise neighbour
        endif
      endif
    endif
    if "leder found" received from counterclockwise neighbour of v then
      if v.leaderStatus ≠ "leader" then
        v.leaderStatus = "not leader"
        send "leder found" to clockwise neighbour
      endif
    endif
  endwhile
endproc
  
```

Ledervalg

```

proc RingLeaderElection
  forall v do
    send v.UID to clockwise neighbour
  endforall

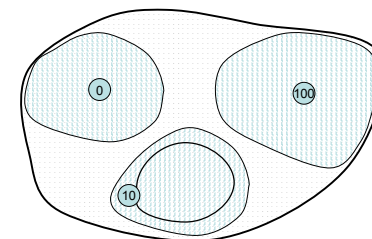
  while not < all v have received "leder found" message > do
    if u.UID received from counterclockwise neighbour of v then
      if v.UID < u.UID then
        send u.UID to clockwise neighbour
      else
        if v.UID = u.UID then
          v.leaderStatus = "leader"
          send "leder found" to clockwise neighbour
        endif
      endif
    endif
    if "leder found" received from counterclockwise neighbour of v then
      if v.leaderStatus ≠ "leader" then
        v.leaderStatus = "not leader"
        send "leder found" to clockwise neighbour
      endif
    endif
  endwhile
endproc
  
```

Tidskompleksitet : $O(n)$
 Meldingskompleksitet : $O(n^2)$

Algoritmer i distribuerte nettverk

I generelle nettverk er det litt mer komplisert enn i ringer...

- Vi kan ikke vite om en node som får tilbake sin egen UID er leder, eller om denne UIDen bare har gått en "liten runde".



- Kjenner vi diameteren d i nettverket, vet vi at vi må vente d steg før vi kan utpeke en leder. (Diameteren til en graf er lengden av lengste korteste sti mellom noe par av noder.)

Ledervalg

```
proc LeaderElection
  forall v do
    v.maxUID = v.UID
  endforall

  for round = 1 to d do
    forall v do
      send v.maxUID to all neighbours in N(v)
      for < each u.maxUID received > do
        if v.maxUID < u.maxUID then
          v.maxUID = u.maxUID
        endif
      endfor
    endforall
  endfor

  forall v do
    if v.maxUID = v.UID then
      v.leaderStatus = "leader"
    else
      v.leaderStatus = "not leader"
    endif
  endforall
endproc
```

Tidskompleksitet : $O(d)$
Meldingskompleksitet : $O(dm)$

Kringkasting og BFS

- Kringkasting av informasjon, og innsamling av informasjon er fundamentale operasjoner innen distribuert og parallell databehandling.
- En node får f.eks ansvar for å les inn data fra en fil, så fordele arbeid rundt på ulike prosessorer, og til slutt samle inn delsvare og gi et svar.

Kringkasting og BFS

Flooding

1. En node som skal kringkaste sender sin melding til alle sine naboer.
 2. Alle noder som mottok meldingen, sender denne videre til sine naboer.
 3. Som sender den videre til sine naboer.
 4. ...
- Brukes i dynamiske nettverk.
 - Ulempen er at det blir mye kommunikasjon.
 - Trenger en metode for å avgjøre når vi er ferdige (f.eks diameter).

Kringkasting og BFS

Kringkasting i trær

I et tre kan vi implementere kringkasting effektivt:

- Roten starter, meldingen sendes nedover i treet, helt til den når bladnodene.
- Hver node må ha en liste over sine barn.

Kringkasting og BFS

Kringkasting i trær

Kjenner vi ikke strukturen i treet (foreldre / barn), kun roten, må vi bygge opp trestrukturen.

1. Bladnodene (untatt en eventuell rot) ser at de bare har *en* nabo *v*, og sier ifra til *v* at de er barn av denne.
2. Når en internnode (untatt en eventuell rot) har motatt barnemeldinger fra alle untatt *en* nabo *w* sier den fra til *w* at den er barn av denne.
3. Vi fortsetter helt til roten har fått barnemeldinger fra alle sine barn.

Kringkasting og BFS

Kringkasting i trær (convergecast)

Beregninger kan nå gjøres bottom-up i treet:

1. Bladnodene sender sin verdi (f.eks en sum) til forelderen.
2. Når en node har fått verdier fra alle sine barn gjøres en beregning, og verdien sendes videre til dennes forelder.

Kringkasting og BFS

Bredde først-trær

Vi har jo som regel ikke tre, men ønsker å legge en trestruktur oppå den generelle grafen (nettverket) vi har – et spennetre.

1. Initielt er treet tomt. Alle noder har en `inTree`-variabel som er FALSE.
2. Roten (utpekt på forhånd) setter `inTree` til TRUE, og sender "join tree" til alle sine naboer.
3. Noder som mottar "join tree", velger ut *en* node den mottok meldingen fra (de kan være flere), og sender "your child" til denne, og "not your child" til de andre den mottok "join tree" fra. `inTree` setter til TRUE, forelderen huskes.
4. I neste runde sender noder som akkurat ble med i treet "join tree" til de av sine barn som enda ikke er med i treet. Er alle naboer med i treet, sendes "terminated" til forelderen. Er "terminated" mottatt fra alle barn, sendes også "terminated" til forelderen.
5. 3 og 4 gjentas til roten mottar "terminated" fra alle sine barn.

Kringkasting og BFS

```
proc BFSDistributed(s)
  s.parent = NIL
  s.inTree = TRUE
  s sends "join tree" to all neighbours

  forall v do
    if v.inTree = FALSE and "join tree" received then
      < choose a node u that has sent "join tree" >
      v.parent = u
      v.inTree = TRUE
      send "your child" to u
      send "not your child" to (N(v) - u)
      send "join tree" to N(v)
    else
      if child status message is received from all neighbours then
        if < v has no children > then
          send "terminated" to v.parent // v er bladnode
        else
          if < "terminated" received from all children > then
            send "terminated" to v.parent
            halt // v stopper
          endif
        endif
      endif
    endif
  endforall
endproc
```

Ledervalg – igjen

Bredde først-trær

Tidligere brukte vi diameter-trikset for å finne en leder i en generell graf. Vi kan bruke spenntreer.

1. For hver node s i grafen, beregner vi et spenntre med s som rot. (Ett tre for hver node, beregnes uavhengig av hverandre.)
2. Hver node s gjøres så til mål for en convergecast av maxUID i treet den er rot i.
3. Noden som ender opp med maxUID lik sin egen UID er leder.

Korteste stier

Bredde først-trær

Et BFS-tre T_s med rot s er et korteste sti-tre (færrest antall kanter) mellom s og de øvrige nodene i grafen.

Vi kan også bruke en distribuert variant av Bellmann-Ford-algoritmen for å finne korteste stier når kantene har ulike vekter.

Korteste stier (Bellmann-Ford)

```
proc BellmannFordDistributed(s)
  forall v do // Feil i boka
    v.dist = ∞
    v.parent = ∞
  endforall

  s.dist = 0
  s.parent = 0

  for round = 1 to n - 1 // n antaes kjent
    forall v do
      send v.dist to all nodes in Nout(v)

      < after receiving u.dist from all u in Nin(v) >
      if v.dist > u.dist + w(uv) for u ∈ Nin(v) then
        v.parent = < u with lowest u.dist + w(uv) > // Utelatt i boka
        v.dist = min {v.dist, u.dist + w(uv)} for u ∈ Nin(v)
      endif
    endforall
  endfor
endproc
```

Korteste stier (Bellmann-Ford)

```
proc BellmannFordDistributed(s)
  forall v do
    v.dist = ∞
    v.parent = ∞
  endforall

  s.dist = 0
  s.parent = 0

  for round = 1 to n - 1 // n antaes kjent
    forall v do
      send v.dist to all nodes in Nout(v)

      < after receiving u.dist from all u in Nin(v) >
      if v.dist > u.dist + w(uv) for u ∈ Nin(v) then
        v.parent = < u with lowest u.dist + w(uv) >
        v.dist = min {v.dist, u.dist + w(uv)} for u ∈ Nin(v)
      endif
    endforall
  endfor
endproc
```

Tidskompleksitet : $O(n)$
Meldingskompleksitet : $O(nm)$

Korteste stier (Floyd)

Korteste stier alle til alle.

Floyds algoritme går i n steg. Nodene har numre 0 til $n-1$, og i k te steg gir vi oss lov til å bruke node k som internnode i stien vi bygger opp mellom to noder i og j .

I den distribuerte varianten antar vi at det finnes en total ordning ($<$) av nodene, f.eks ved hjelp av UID. (Prosessorer i et datanett har alltid en slags ID, f.eks en IP-adresse. Helt anonyme nett mer en slags matematisk finurlighet.)

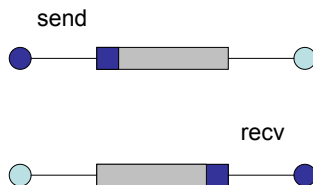
Korteste stier (Floyd)

```
proc FloydDistributed()
  forall v do
    forall u do
      if u ∈ Nout then
        v.distTable[u] = w(vu)
        v.firstTable[u] = u
      else
        if u = v then
          v.distTable[u] = 0
        else
          v.distTable[u] = ∞
        endif
      endif
    endforall
  endforall

  for pass = 0 to n - 1 do
    < let p be the next node in the node ordering >
    < p broadcasts p.distTable to all other nodes >
    forall v do
      forall u do
        if v.distTable[u] > v.distTable[p] + p.distTable[u] then
          v.distTable[u] = v.distTable[p] + p.distTable[u]
          v.firstTable[u] = v.firstTable[p]
        endif
      endforall
    endforall
  endfor
endproc
```

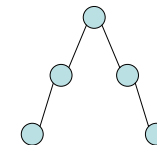
Asynkron modell

- Vi har hittil antatt at algoritmene våre går stegvis. Innenfor hvert steg har det ikke vært så nøye hvor fort ting gjøres, så lenge vi venter til alle er ferdige med å gå til neste steg.
- Dette er kanskje ikke den mest nøyaktige modellen av Den Virkelige Verden.
- Vi lar kommunikasjonskanalene i nettverket (kantene i grafen) ha buffere for meldinger.



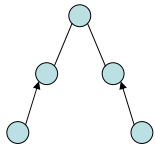
Asynkront ledervalg i trær

Algoritmen er essensielt den samme som den synkrone, med unntak av at kommunikasjon implementeres med buffere, og en litt annen slutttilstand.



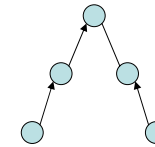
Asynkront ledervalg i trær

Algoritmen er essensielt den samme som den synkrone, med unntak av at kommunikasjon implementeres med buffere, og en litt annen slutttilstand.



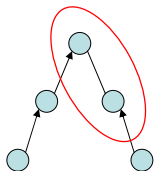
Asynkront ledervalg i trær

Algoritmen er essensielt den samme som den synkrone, med unntak av at kommunikasjon implementeres med buffere, og en litt annen slutttilstand.



Asynkront ledervalg i trær

Algoritmen er essensielt den samme som den synkrone, med unntak av at kommunikasjon implementeres med buffere, og en litt annen slutttilstand.



Disse to må bli enige seg imellom.
I praksis lar man noden med lavest
(eller høyest) UID bli leder.

Asynkront ledervalg i generelle grafer

- Vi kan ikke lengre bruke diameter-trikset for å avgjøre om algoritmen er ferdig, ettersom den ikke lengre går i pent definerte steg.
- Vi kan imidlertid bruke spenntre-metoden, og ledervalg i trær.

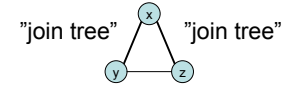
Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



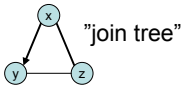
Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



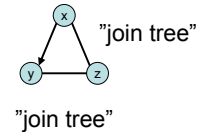
Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



Ikke et korteste sti-tre.

Asynkron kringkasting

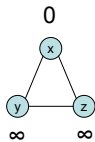
Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.

Modifikasjonen består i å vedlikeholde en avstandsvariabel $d(v)$ for hver node v .

1. Rotnoden starter ved å sende avstand 0 til sine naboer.
2. Hver node v ser så på meldingene den mottar fra sine naboer. For hver melding hentet fra uv sitt sendBuffer, "relaxer" vi kanten uv . Dvs. at hvis $d(v) > d(u) + 1$, så setter vi $d(v) = d(u) + 1$, og $\text{parent}(v) = u$.
3. Etter å ha scannet bufferne til naboene og relaxet kantene, sender v oppdatert informasjon til sine naboer.

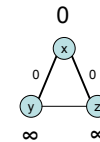
Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



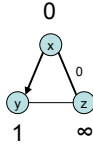
Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



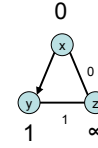
Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



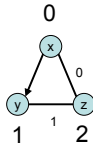
Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.



Asynkron kringkasting

Algoritmen for å finne bredde først-trær i grafer kan ikke benyttes uten modifikasjoner.

