

Election algorithms

- Definition of election algorithm, assumptions
- election algorithm on a tree
 - LeLann
 - Chang and Roberts
 - Dolev-Klawe-Rodeh
- arbitrary networks - extinction algorithm

1

What is election algorithm

- Election algorithm satisfies the following properties:
 - (uniformity) each process has the same local algorithm
 - (decentralization) a computation can be initiated by arbitrary subset of processes
 - (termination) each computation terminates
 - (safety) in every terminal state there is one and only one process in the state of leader
- Assumptions
 - system is fully asynchronous - not necessary but useful
 - each process has unique identity (name), process identities are totally ordered (can be compared) - needed to break the symmetry, smallest identity wins (becomes leader)
 - the only operation allowed on identities is comparison - in asynchronous systems arbitrary (non-comparison algorithms can do no better)
 - each message may contain up to a constant number of identities - so complexities of algorithms can be compared

2

Waves on trees and election

- Use wave tree algorithm
- first phase - "wake up" all processes (non-initiators)
 - a wakeup message is propagated from initiators to other processes
- second phase - do a wave on the tree
 - leaves start
 - id of lowest process in subtree is attached to token
- message complexity: $4N-4$
 $O(N)$ (two wakeup messages and two tokens are sent along each channel)
- time complexity - $3D+1$
 $(D - \text{diameter}) - O(N)$
 - D - to send wakeup messages, in $D+1$ tree algorithm starts
 - D - to make first decision, D - to propagate decision

```

var  $wsp$  : boolean      init false ;
 $wrp$  : integer          init 0 ;
 $recp[q]$  : boolean for each  $q \in Neigh_p$  init false ;
 $v_p$  :  $\mathcal{P}$                init  $p$  ;
 $state_p$  : (sleep, leader, lost)      init sleep ;

begin if  $p$  is initiator then
  begin  $wsp := true$  ;
    forall  $q \in Neigh_p$  do send (wakeup) to  $q$ 
    end ;
  while  $wrp < \#Neigh_p$  do
    begin receive (wakeup) ;  $wrp := wrp + 1$  ;
      if not  $wsp$  then
        begin  $wsp := true$  ;
          forall  $q \in Neigh_p$  do send (wakeup) to  $q$ 
          end
        end ;
      (* Now start the tree algorithm *)
      while  $\#(q : \neg recp[q]) > 1$  do
        begin receive (tok,  $r$ ) from  $q$  ;  $recp[q] := true$  ;
           $v_p := \min(v_p, r)$ 
        end ;
        send (tok,  $v_p$ ) to  $q_0$  with  $\neg recp[q_0]$  ;
        receive (tok,  $r$ ) from  $q_0$  ;
         $v_p := \min(v_p, r)$  ; (* decide with answer  $v_p$  *)
        if  $v_p = p$  then  $state_p := leader$  else  $state_p := lost$  ;
        forall  $q \in Neigh_p$ ,  $q \neq q_0$  do send (tok,  $v_p$ ) to  $q$ 
        end
      end
    end
  end
end

```

3

LeLann's algorithm

```

var  $List_p$  : set of  $\mathcal{P}$       init { $p$ } ;
 $state_p$  ;

begin if  $p$  is initiator then
  begin  $state_p := cand$  ; send (tok,  $p$ ) to  $Next_p$  ; receive (tok,  $q$ ) ;
    while  $q \neq p$  do
      begin  $List_p := List_p \cup \{q\}$  ;
        send (tok,  $q$ ) to  $Next_p$  ; receive (tok,  $q$ )
      end ;
      if  $p = \min(List_p)$  then  $state_p := leader$ 
      else  $state_p := lost$ 
    end
    else repeat receive (tok,  $q$ ) ; send (tok,  $q$ ) to  $Next_p$  ;
      if  $state_p = sleep$  then  $state_p := lost$ 
      until false
    end
  end
end

```

- Each initiator sends token with its id around the ring
- a process forwards foreign token and records the identity it carries
- a process decides when it receives its own token
- if processes gets a message before it sends its own - it loses
- complexity
 - message - $O(N^2)$
 - time - $O(N)$

4

Chang and Roberts's algorithm (CR)

- Process propagates foreign message only when the original sender may become a leader
- complexity
 - time - same as LeLann's
 - message
 - worst case - same as LeLann's
 - average - $O(N \log N)$

```

var  $state_p$  ;

begin if  $p$  is initiator then
  begin  $state_p := cand$  ; send (tok,  $p$ ) to  $Next_p$  ;
    repeat receive (tok,  $q$ ) ;
      if  $q = p$  then  $state_p := leader$ 
      else if  $q < p$  then
        begin if  $state_p = cand$  then  $state_p := lost$  ;
          send (tok,  $q$ ) to  $Next_p$ 
        end
      end
    until  $state_p = leader$ 
  end
  else repeat receive (tok,  $q$ ) ; send (tok,  $q$ ) to  $Next_p$  ;
    if  $state_p = sleep$  then  $state_p := lost$ 
    until false
  end
end
(* Only the leader terminates the program. It floods a message to all processes to inform them of the leader's identity and to terminate *)

```

5

Dolev-Klawe-Rodeh algorithm (DKR)

- Worst case message complexity of Chang-Roberts is still $O(N^2)$ -- in the worst case the identities that are not leaders are allowed to propagate
- idea of DKR - eliminate wrong identities as soon as possible
- active/passive processes as in CR
 - active - propagates its current identity (stored in variable c_i) to downstream processes
 - passive - forwards messages
- algorithm proceeds in rounds, each round has two phases:
 - propagation - each active process sends message $\langle one, c_i \rangle$ to downstream active process, received id stored in acn
 - elimination - each active process sends a message $\langle two, acn \rangle$ to downstream neighbor when $\langle two, acn \rangle$ arrives, the receiver compares it with its can thus the following identity "catches up" with preceding
- when process gets his own identity - its is a winner
 winner sends message $\langle small, id \rangle$ informing others

6

DKR

```

var  $c_i$  :  $\mathcal{P}$  init  $p$ ; (* Current identity of  $p$  *)
 $acn_p$  :  $\mathcal{P}$  init undef; (* Id of anticlockwise active neighbor *)
 $win_p$  :  $\mathcal{P}$  init undef; (* Id of winner *)
 $state_p$  (active, passive, leader, lost) init active;

begin if  $p$  is initiator then  $state_p := active$  else  $state_p := passive$ ;
  while  $win_p = undef$  do
    begin if  $state_p = active$  then
      begin send (one,  $c_i$ ); receive (one,  $q$ );  $acn_p := q$ ;
        if  $acn_p = c_i$  then (*  $acn_p$  is the minimum *)
          begin send (small,  $acn_p$ ); receive (small,  $q$ );
            then  $c_i := acn_p$ ;
          end
        else (*  $acn_p$  is current id of neighbor *)
          begin send (two,  $acn_p$ ); receive (two,  $q$ );
            if  $acn_p < c_i$  and  $acn_p < q$ 
              then  $c_i := acn_p$ 
              else  $state_p := passive$ 
            end
          end
        end
      end
    else (*  $state_p = passive$  *)
      begin receive (one,  $q$ ); send (one,  $q$ );
        receive  $m$ ; send  $m$ ;
        (*  $m$  is either (two,  $q$ ) or (small,  $q$ ) *)
        if  $m$  is a (small,  $q$ ) message then  $win_p := q$ 
        end
      end
    end
  end
  if  $p = win_p$  then  $state_p := leader$  else  $state_p := lost$ 
end

```

DKR message complexity estimate

- Claim: if there were k identities at the beginning of the round
 - only $k/2$ survive after the round
 - ♦ each process (if it survives round) assumes the identity of the first upstream neighbor
 - ♦ out of two neighbor identities only one survives the round
- there can be at most $\log N$ rounds, $2N$ messages are exchanged during round, algorithm completes with N messages informing processes of results
 - ♦ thus message complexity of the algorithm is:

$$2N \log N + N = O(N \log N)$$

8

Election on arbitrary networks, extinction

- Election can be done on networks of arbitrary topology by using any centralized wave algorithm.
- The technique is called *extinction*.
- each initiator starts a separate wave. To distinguish waves, initiator's id is attached to token.
- Only the wave of the initiator with the lowest id is allowed to finish
- A process stores the wave with the lowest id that it saw in currently active wave (*caw*) variable. The process ignores tokens from waves with higher id than *caw*
- If a process receives a token from a wave with still lower id than *caw*, the process copies that id into *caw*
- complexity
 - ♦ message NM where M - number of messages used in the base wave algorithm
 - ♦ time NT where T - time the wave takes

9

Extinction applied to Echo

```

var  $caw_p$  :  $\mathcal{P}$  init undef; (* Currently active wave *)
 $rec_p$  : integer init 0; (* Number of (tok,  $caw_p$ ) received *)
 $father_p$  :  $\mathcal{P}$  init undef; (* Father in wave  $caw_p$  *)
 $lrec_p$  : integer init 0; (* Number of (ldr, ...) received *)
 $win_p$  :  $\mathcal{P}$  init undef; (* Identity of leader *)

begin if  $p$  is initiator then
  begin  $caw_p := p$ ;
    for all  $q \in Neigh_p$  do send (tok,  $p$ ) to  $q$ 
    end
  while  $lrec_p < \#Neigh_p$  do
    begin receive msg from  $q$ ;
      if msg = (ldr,  $r$ ) then
        begin if  $lrec_p = 0$  then
          for all  $q \in Neigh_p$  do send (ldr,  $r$ ) to  $q$ ;
             $lrec_p := lrec_p + 1$ ;  $win_p := r$ 
          end
        else (* a (tok,  $r$ ) message *)
          begin if  $r < caw_p$  then (* Reinitialize algorithm *)
            begin  $caw_p := r$ ;  $rec_p := 0$ ;  $father_p := q$ ;
              for all  $s \in Neigh_p$ ,  $s \neq q$ 
                do send (tok,  $r$ ) to  $s$ 
              end
            end
          if  $r = caw_p$  then
            begin  $rec_p := rec_p + 1$ ;
              if  $rec_p = \#Neigh_p$  then
                if  $caw_p = p$ 
                  then for all  $s \in Neigh_p$ 
                    do send (ldr,  $p$ ) to  $s$ 
                  else send (tok,  $caw_p$ ) to  $father_p$ 
                end
                (* If  $r > caw_p$ , the message is ignored *)
              end
            end
          end
        end
      end
    end
  end
  if  $win_p = p$  then  $state_p := leader$  else  $state_p := lost$ 
end

```