

A Distributed Spanning Tree Algorithm

Karl Erik Johansen
Ulla Lundin Jørgensen
Svend Hauge Nielsen
Søren Erik Nielsen
Sven Skyum

DAIMI PB – 226
Marts 1987

A Distributed Spanning Tree Algorithm

Abstract

We present a distributed algorithm for constructing a spanning tree for connected undirected graphs. Nodes correspond to processors and edges correspond to two way channels. Each processor has initially a distinct identity and all processors perform the same algorithm. Computation as well as communication is asynchronous. The total number of messages sent during a construction of a spanning tree is at most $2E+3N\log N$. The maximal message size is $\log\log N + \log(\text{maxid}) + 3$, where maxid is the maximal processor identity.

1. Introduction

Construction of spanning trees for communication graphs has proven useful and has been considered in a number of papers ([3],[4],[5],[6],[7]). Other problems such as termination, extremafinding, and election of a leader all reduce to spanning trees. Most papers on spanning trees deals with construction of minimum spanning trees. Santoro has shown that $O(E+N\log N)$ is a lowerbound on the message-complexity for the problem. The best known upper bound is $2E+5N\log N$ ([3]). Recently Lavee and Roucairol have presented an algorithm that constructs a spanning tree in a general network with message complexity $N-1+3N\log N$ provided the algorithm behaves in a balanced manner ([6]). Their worst-case complexity is again $2E+5N\log N$ and their message-size is very large. Finally Korach, Moran and Zaks have shown that finding a spanning tree in a complete graph might be easier than finding a minimal spanning tree ([5]). In this paper we present an algorithm, having worst-case complexity $2E+3N\log N$, that constructs a spanning tree in an arbitrary network of processors.

The algorithm is based on the commonly used model. We consider an undirected connected graph without selfloops. Each node corresponds to a processor with unique identity, and each edge corresponds to a two way channel. Each channel has (input) buffers at either endpoint organized as queues. Processors can send and receive messages via channels. The communication is asynchronous. Messages on a channel are received at an input buffer for a processor in the order they are sent from the neighbour, they may be arbitrarily but finitely delayed. Initially all processors are in a sleeping state. They wake up spontaneously after a finite time and start execution of the algorithm. The algorithm is based on a finite state machine where states have a little memory. In the papers referenced above processors can either wake up spontaneously or be waked up by receiving messages. This difference in presentation makes no "visible" difference in the behavior, since messages can be

delayed arbitrarily.

Section 2 contains various concepts and notions needed to explain the algorithm. In Section 3 an overall description of the algorithm is given. The algorithm is given in details in Section 4. Section 5 contains the proof of correctness while Section 6 contains the analysis of the algorithm.

2. Definitions and notations

Let $G=(V(G),E(G))$ be an undirected connected graph without selfloops. We will use $\{.,.\}$ to denote undirected edges and $(.,.)$ to denote directed edges. Each node v in G corresponds to a processor with the unique identity $id_v=v$. $V(G) = \{1,2, \dots\}$. Let $N=|V(G)|$ and $E=|E(G)|$.

A **fragment** is a connected subgraph, $F=(V(F),E(F))$ of G , with no cycles (an undirected tree). A **spanning tree** for G is a fragment F , where $V(F)=V(G)$. A **spanning forrest** is a set (F_1,F_2,\dots,F_k) of fragments, such that $V(F_i) \cap V(F_j)=\emptyset$ for $i \neq j$ and $V(G)=V(F_1) \cup V(F_2) \cup \dots \cup V(F_k)$. Given a fragment $F=(V(F),E(F))$ and a node v in $V(F)$ let $F_v=(V(F),D(F_v))$ denote the rooted tree where v is the root, and edges in $D(F_v)$ are directed towards the root v . Each fragment F is equipped with two not necessarily different orientations by naming two roots. One is called the **centre** of F and denoted $c(F)$. If (v,w) is an edge in $D(F_{c(F)})$ then $\{v,w\}$ in $E(F)$ will be called an **in-edge for v** and an **out-edge for w** . The other is called the **king** of F and denoted $k(F)$. If (v,w) is an edge in $D(F_{k(F)})$ then we call $\{v,w\}$ an **up-edge for v** and a **down-edge for w** . Each fragment F has a unique identification $\langle level,k(F) \rangle$ called the **colour** of F . level will be an integer in $[0,\log N]$.

3. Description of the algorithm

Initially each node v in the network is in a **sleeping** state. It wakes up spontaneously and enters after initialization an **idle** state. Then it constitutes a fragment of size 1 at level 0. The colour is $\langle 0,v \rangle$ and v is both the centre and the king. During execution of the algorithm each centre in the idle state attempts to send a **request(colour)**-message along one of its unprocessed edges in order to combine fragments into larger fragments. If a centre has no adjacent unprocessed edge, the centre is moved around in a fragment F in a depth-first fashion in $F_{k(F)}$ by sending **movecentre**-messages until an unprocessed adjacent edge is found or the algorithm terminates. Assume that we have two fragments F_1 and F_2 with colours $\langle L_1,k(F_1) \rangle$ and $\langle L_2,k(F_2) \rangle$. Colours of nodes in F_i are then at most, but not necessarily equal to $\langle L_i,k(F_i) \rangle$. During computation they will all receive the colour $\langle L_i,k(F_i) \rangle$ sooner or later. Assume furthermore that the centre c_1 in F_1 sends a request($\langle L_1,k(F_1) \rangle$)

to a node in F_2 along edge e (see Figure 1). After sending the request c_1 enters state **waiting_for_accept**. The request is routed in the direction of the centre c_2 in F_2 as long as $\langle L_1, k(F_1) \rangle$ is greater than the colour of the nodes it passes. Nodes which pass

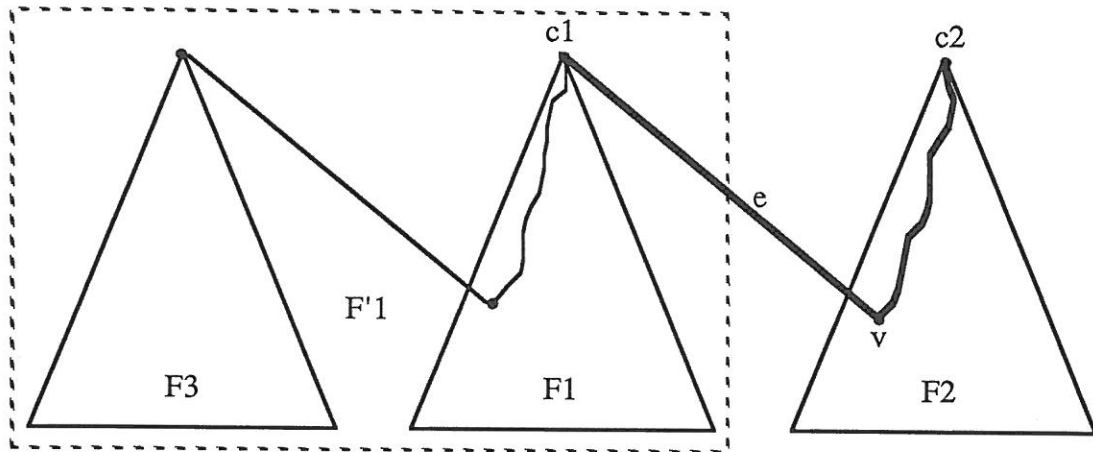


Figure 1. Connecting fragments.

on request-messages also enter the **waiting_for_accept** state and then only listen to the in-edge for an **accept-** or **newcolour-**message (see later). If a request arrives at a node with greater colour, the node does not react on the request. In that case nodes which passed on that request-message will receive a new colour and reenter the idle state. If the request reaches the centre c_2 , call the route (in both directions) along which the request came for the **request-route**. If $\langle L_2, k(F_2) \rangle$ is less than $\langle L_1, k(F_1) \rangle$ then F_2 and what has become of F_1 (see later) are to be combined into a larger fragment F . To perform a combination c_2 sends an **accept(L_2)**-message back through F_2 along the request-route. Nodes on the route, including c_2 itself, enter state **waiting_for_new_colour**. When c_1 receives the **accept**-message it might have become centre in a larger fragment, than it was when it sent the request to F_2 . It might also be in state **waiting_for_new_colour** in which case F_1 is in the process of being combined with yet another fragment F_3 and will receive a new colour in connection with that combination. In both cases c_1 has sent an **accept**-message after sending the request-message to F_2 . This possibility is necessary to prevent dead-locks. If c_1 is waiting for a new colour, it waits until it has received a new colour before it goes on, otherwise it immediately initiates the final part of the combination of F'_1 and F_2 , where F'_1 is the present fragment containing c_1 as its centre. The combination of F'_1 and F_2 will be a fragment F consisting of F'_1 , F_2 and the edge e along which c_1 sent the request to F_2 . The centre for F will be c_2 , the king will be $k(F'_1)$ and the new level L will be L'_1 , if L'_1 is greater than L_2 , and L'_1+1 otherwise. (L'_1 is the level of F'_1). c_1 stops being a centre and starts colouring F_2 by sending a **newcolour($\langle L, k(F'_1) \rangle$)**-message along e . The **newcolour**-message is

broadcast to all nodes in F_2 via tree-edges. During colouring of F_2 the orientation with respect to up and down in $F_k(F)$ is updated. When c_2 receives its new colour, it reenters the idle state and we say that F has been formed. If $\langle L, k(F_1) \rangle$ differs from c_1 's colour, F_1 is also coloured by sending newcolour-messages. The reason for choosing $k(F_1)$ as king for F is that nodes in F_1 do not necessarily receive a new colour in connection with the combination.

If a centre c in a fragment F sends a request along an unprocessed edge e to a node v in the same fragment, then v 's colour might be less than the colour sent with the request and v cannot know that it comes from the same fragment (see Figure 2). The request might therefore be forwarded in the direction of the centre as described above. Sooner or later the request-message will meet a node u with the same colour as c (it might be c itself) and will not be sent any further. When v later on receives the colour of c , it recognizes that the request it forwarded was received from the centre of its own fragment and it sends a **close**-message back along the edge e to c and marks the edge closed. Upon reception of the close-message c marks the edge e closed as well and reenters the idle state.

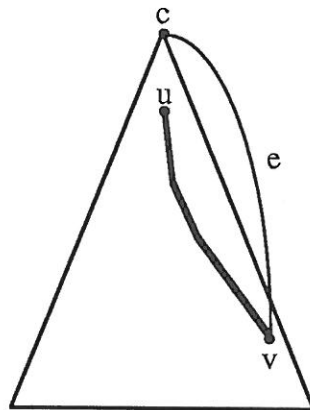


Figure 2. The route followed by a request from a centre to a node in its own fragment

The preceding description gives the overall picture of the algorithm but due to parallelism, messages might cross each other, which adds tedious details to the algorithm.

4. The algorithm

During computation, nodes (or processors) mark their adjacent edges (or channels) with attributes. More attributes can be attached to each edge and the attributes at the two endpoints might differ. The possible attributes are:

open. All edges adjacent to a node are marked open by that node when it wakes up. Messages will only be sent along open edges. (Edges which are only marked open have been and will be referred to as unprocessed).

closed. No messages will be sent along closed edges.

branch. A node marking an edge, with the attribute branch, knows that the edge is part of the fragment and will be part of the final spanning tree. (Edges marked branch at both ends have been and will be referred to as tree-edges).

in, out, up, and down. The edge is an in-edge (out-, up-, or down-edge resp. See Section 2). Only branch-edges will be attached those attributes.

Nodes can send a number of different messages along open edges. There are five different message types that can be send, namely:

request, accept, close, newcolour, and movecentre. Request and newcolour carry a colour as parameter, accept has a level as parameter, while close and movecentre are without parameters.

Nodes can be in a number of different states $\langle st_1, st_2 \rangle$. Apart from the various attributes attached to edges mentioned above, only information about at most one edge and one colour has to be stored in a node. We have chosen to store that information within the states.

st_1 is either **centre** or **ordinary** while st_2 is one of the following:

sleeping. Initially all nodes are sleeping.

waiting_for_accept(colour,e). The node is a centre and has sent a request(colour) along e or it is not a centre and has forwarded a request(colour) which was received along e.

waiting_for_new_colour(e). The node has sent (or forwarded) an accept-message along e and waits for a new colour.

terminated. The node knows all its adjacent tree-edges and has finished its participation in the distributed computation.

idle. The node takes part in the computation, but is in none of the former four states.

At termination of the distributed algorithm (all nodes are in state terminated) all edges will be marked closed at both ends. Edges will be either tree-edges or not being marked branch in either endpoint in which case they are called cross-edges.

Each node or processor v in state $\langle st_1, st_2 \rangle$ executes the following algorithm. The fragment including v is referred to by F and has colour= $\langle L, id \rangle$. $A(v)$ denotes the adjacent edges. The edges in $A(v)$ are ordered such that, if we choose an edge with a specific property, it is assumed that we always choose the first edge in the ordering with that property. This is important when requests are sent from idle centres.

```

repeat
case <st1,st2> of
<*,sleeping> :
  <st1,st2>:= <centre,idle>; colour:=<0,v>;
  for all e in A(v) do mark(e)={open} od;

<centre,idle> , <centre,waiting_for_accept(colour,e)> :
  if inputbuffer for an edge e1 in A(v) is nonempty then m:=read buffer e1;
  case m of
  close: {if st2 = waiting_for_accept then e1=e}
    mark(e1)=mark(e1)-{open}+{closed}; st2:=idle;
  request(colour1):
    if colour < colour1 then send accept(L) along e1;
    st2:=waiting_for_new_colour(e1) fi;
  accept(L1): {if st2 = waiting_for_accept then e1=e}
    if L1 = L then L:=L+1; colour:=<L,id>;
    for all e2 in A(v) where out in mark(e2) do
      send newcolour(colour) along e2 od fi;
    send newcolour(<L,id>) along e; mark(e)={open,down,in,branch};
    <st1,st2>:=<ordinary,idle>;
  endcase
  else if st2=idle then
    if there is an e in A(v) where mark(e)={open} then
      send request(colour) along e; st2:=waiting_for_accept(colour,e)
    else if there is an e in A(v) where {open,down} < mark(e) then
      send movecentre along e; mark(e)=mark(e)-{out}+{in};
      st1:=ordinary
    else if there is an e in A(v) where {open,up} < mark(e) then
      send movecentre along e;
      mark(e)=mark(e)-{open,out}+{closed,in};
      <st1,st2>:=<ordinary,terminated>
    else st2:=terminated fi fi fi fi;

<*,waiting_for_new_colour(e)>:
  if inputbuffer for e is nonempty then m:=read buffer e;
  case m of
  newcolour(colour1):
    colour:=colour1; mark(e)={open,up,out,branch};
    for all e1≠e in A(v) where branch in mark e1 do
      send newcolour(colour) along e1;
      if up in mark(e1) then mark(e1)=mark(e1)-{up}+{down} fi od;
    st2:=idle;
  request(colour1): {skip};
  endcase fi;

```



```

<ordinary,idle> :
  if an inputbuffer for an edge e in A(v) is nonempty then m:=read buffer e;
  case m of
  request(colour1) :
    if colour < colour1 then
      for the in-edge e1 in A(v) do
        if e≠e1 then send request(colour1) along e1;
          st2:=waiting_for_accept(colour1,e) fi od
        else if (colour = colour1) and mark(e) = {open} then
          send close along e; mark(e):= {closed} fi fi;
newcolour(<L1,id1>) :
  colour:=<L1,id1>;
  if (down in mark(e)) and (id1≠id) then mark(e):=mark(e)-{down}+{up}
  fi;
  for all e1≠e in A(v) where branch in mark(e) do
    send newcolour(colour) along e1;
    if (up in mark(e1)) and (id1≠id) then
      mark(e1):=mark(e1)-{up}+{down} fi od;
movecentre :
  if down in mark(e) then mark(e):=mark(e)-{open}+{closed} fi;
  mark(e):=mark(e)-{in}+{out}; st1:=centre
endcase fi;

```

```

<ordinary,waiting_for_accept(colour1,e)> :
  if the inputbuffer for the in-edge e1 in A(v) is nonempty then
  m:=read buffer e1;
  case m of
  request(colour2): {This might occur if v has just been a centre - skip}
  accept(L) :
    send accept(L) along e; st2:=waiting_for_new_colour(e)
  newcolour(colour2) :
    if colour1 <= colour2 then st2:=idle;
      if (colour1= colour2) and mark(e)={open} then
        send close along e; mark(e):={closed} fi fi;
      for all out-edges e2 in A(v) do send newcolour(colour2) along e2;
        if (up in mark(e2)) and (id2≠id) then
          mark(e2):=mark(e2)-{up}+{down} fi od;
    colour:=colour2;
  movecentre:
    if down in mark(e1) then mark(e):=mark(e)-{open}+{closed} fi;
    mark(e):=mark(e)-{in}+{out}; send accept(L) along e;
    <st1,st2>:=<centre,waiting_for_new_colour(e)>;
  endcase fi;
endcase until st2=terminated;

```

5. Correctness of the algorithm

For each node at most one message is read during an execution of a cycle of the algorithm. We may therefore w.l.o.g. assume that time is discrete $(-N, \dots, -1, 0, 1, 2, \dots)$ and exactly one node executes one cycle of the algorithm for each time instance t . We may furthermore assume that all nodes are awake and that no messages have been read at time 0.

Lemma 5.1

When a node v terminates then all nodes u , which can be reached from v following down-edges (coincide with out-edges), will be terminated as well (and have no open adjacent edges).

Proof

For $t=0$ no node is terminated, so the Lemma trivially holds true. If a node v terminates at time t , then v has at most one open adjacent edge (an up-edge). Since all down-edges $\{v, w\}$ are closed v has received a movecentre-message along these edges indicating that "down-neighbours" w are terminated. Thus the Lemma follows by induction in t .

Corollary 5.2

If a node u has an open adjacent edge, a nonterminated centre c (possibly u itself) can be reached from u following in-edges marked open.

Lemma 5.3

If there is more than one centre at time t , then every pair of centres c_1 and c_2 are connected by an open path (a path where all edges are marked open at both endpoints).

Proof

The Lemma holds true for $t=0$ because the network is connected and all edges are open. The Lemma will remain true from time t to $t+1$ if the node v executing its cycle at time t does not mark any new edge closed. Therefore assume that the Lemma holds true at time t and that node v closes an edge at time t .

There are four possibilities: (1) v closes an up-edge and terminates, (2) v closes a down-edge $\{v, w\}$ after receiving a movecentre-message along $\{v, w\}$, (3) v closes $\{v, w\}$ if w is centre in the fragment containing v at time t and v has received a request from w with v 's colour or v has received a newcolour-message with the same colour as the colour of an earlier request from w , or (4) v closes $\{v, w\}$ after receiving a close-message along $\{v, w\}$.

Ad (1): By Lemma 5.1 no path of open edges goes through v so the connectivity of

centres with respect to open edges is not affected by this operation.

Ad (2): Receiving a movecentre-message along a down-edge $\{v,w\}$ indicates that w has terminated and again by Lemma 5.1 we get that $\{v,w\}$ does not contribute to the connectivity of centres.

Ad (3): Before v closes $\{v,w\}$ at time t , there exists a cycle of open edges containing v and w . Breaking this cycle does not affect the connectivity either.

Ad (4): Similar to case 3.

Lemma 5.4

No cycle of tree-edges can be formed.

Proof

A cycle could only be formed if a centre c would accept a request that was initiated by c itself. That will never happen since the colour of a node is nondecreasing during computation.

Lemma 5.5

The number of tree-edges equals the difference between the total number of nodes and the number of centres.

Proof

The Lemma holds true initially and creation of a new tree-edge and deletion of a centre happen at the same time instance for a node in state $\langle \text{centre}, \text{waiting_for_accept}(*,*) \rangle$ after receiving an accept-message.

Theorem 5.6

The algorithm will terminate (all nodes are terminated). At termination all edges are closed and the tree-edges form a spanning tree.

Proof

The analysis of the number of messages sent (see Section 6) implies that the network will reach a stable situation, where no more messages will be sent and no more computing go on. Let the network be stable at time t . Assume that c is a nonterminated centre of maximal colour present in the network at time t , if such one exists. c cannot be idle because an idle centre can execute a cycle of computation in all circumstances. If c is waiting for an accept along $\{c,v\}$, then $\{c,v\}$ would be open and an open path from v to a nonterminated centre c_1 (possibly c) along in-edges would exist (Corollary 5.2). The maximality of c 's colour then implies that an accept- or close-message will be sent along $\{c,v\}$ to c at a later time than t , which is a contradiction. If c is waiting for a new colour, it will eventually receive one, so this is impossible as well. Thus all present centres will be terminated and by Corollary 5.2 all edges

will be closed. Lemma 5.4 then implies that exactly one centre exists. It finally follows from Lemmas 5.3 and 5.5 that the set of tree-edges form a spanning tree.

Remark. At the time of termination one node knows that the algorithm is terminated, namely the node terminating into state $\langle \text{centre}, \text{terminated} \rangle$.

6. Analysis of the algorithm

The first Lemma follows directly from the way levels are computed.

Lemma 6.1

If a centre node c has level k , then the number of nodes u (including c itself), which can be reached from c following out-edges, is at least 2^k .

Corollary 6.2

The maximal level obtained during a computation is bounded by $\log N$.

Lemma 6.3

The number of times a node v can receive a newcolour-message is bounded by $\log N$, so the total number of newcolour-messages sent during a computation is bounded by $N \cdot \log N$.

Proof

Follows directly from Corollary 6.2.

Lemma 6.4

The number of close-messages sent during a computation is $E - N + 1$.

Proof

Exactly one close-message is sent for each cross-edge.

Lemma 6.5

The total number of request-messages sent during a computation is bounded by $N \cdot \log N + E - N + 1$.

Proof

A node cannot send two requests without having received a newcolour- or close-message between them.

Lemma 6.6

The total number of movecentre-messages sent along up-edges during a computation is $N-1$.

Proof

All nodes but the centre in the final spanning tree send exactly one such message.

Accept- and movecentre-messages are only sent along edges which become tree-edges. As for tree-edges we are able to bound the total number of accept-messages and movecentre-messages sent along down-edges (**mc**-messages).

Lemma 6.7

For all edges $e=\{v,w\}$ which become tree-edges the total number of accept-messages sent along e and movecentre-messages sent along e in a direction marked down is bounded by $\log N$. Thus the total number of accept-messages and movecentre-messages sent along down-edges (**mc**-messages) during a computation is bounded by $(N-1) \cdot \log N$.

Proof

Let $e=\{v,w\}$ be an arbitrary edge in G along which at least one accept-message is sent. It might happen that both a **mc**- and an accept-message are sent along e between two consecutive newcolour-messages, but we will prove that after sending an accept-message along e , at least two newcolour-messages will be sent before a **mc**-message can be sent. Since two accept-messages or two **mc**-messages clearly cannot be sent between two consecutive newcolour-messages, the Lemma will follow.

Now assume that at time t an accept-message is sent from v to w along e . w is then in state `waiting_for_accept(*,e)` and v enters state `waiting_for_new_colour(e)`. The next message sent along e is therefore a newcolour-message sent from w to v . w enters the idle state and marks e `{branch,out,up}`. On reception of the newcolour-message v enters state `idle` and marks e `{branch,in,down}`. That is, the centre is situated at the "w side" of e and the king is at the "v side" of e . Before a **mc**-message can be sent along e (from w to v) the up and down markings of e have to be changed. This can only happen in connection with another newcolour-message sent along e (from w to v , the in and out markings on e can only be changed in connection with movecentre-messages).

Theorem 6.8

The total number of messages sent during a computation is bounded by $2E+3N \cdot \log N - N - \log N + 1$.

Proof

Lemmas 6.1 through 6.7.

Remark. It is possible to cut $c \cdot N$ off the bound for some $c \geq 1$ by a more careful analysis.

Theorem 6.9

The total number of bits sent during a computation is bound by $(E+2N \cdot \log N) \cdot (\log(\text{maxid}) + \log \log N + 6)$, where maxid is the maximal identification of a node.

Proof

Three bits suffice to give the type of message. The information sent by request- and newcolour-messages is a colour, which is at most $\log(\text{maxid}) + \log \log N$ bits. Accept-messages have a level attached to them. By observing that information on the level is superfluous, except when the accept is sent from one segment to another and the level to be sent in that case is identical to the level of the node sending the message, we get that only

$N-1$ of the accept-messages require $\log \log N$ bits.

References

- [1] Bodlaender H.L., van Leeuwen J.: New Upperbound for Decentralized Extrema Finding in a Ring of Processors. RUU-CS-85-15 Univ. of Utrecht. 1985. (Preliminary version)
- [2] Francez N.: Distributed termination. ACM Trans. on Programming Languages and Systems. Vol. 2. 1980, Pages 42-55.
- [3]. Gallager R. G., Humblet P. A., Spira P. M.: A Distributed Algorithm for Minimum-Weight Spanning Trees. ACM Trans. on Programming Languages and Systems. Vol. 5. 1983, Pages 66-77.
- [4] Humblet P. A.: A Distributed Algorithm for Minimum Weight Directed Spanning Trees. IEEE Trans. on Communications. Vol 31. 1983, Pages 756-762.
- [5] Korach E., Moran S., Zaks S.: The Optimality of Distributive Constructions of Minimum Weight and Degree Restricted Spanning Trees in a Complete Network of Processors. Proc. of 4. Ann. ACM SPDC, Ontario 1985, Pages 277-286.

- [6] Lavellee I., Roucairol G.: A Fully Distributed (Minimal) Spanning Tree Algorithm. *Inf. Proc. Letters* 23. 1986, Pages 55-62.
- [7] Santoro N.: On the Message Complexity of Distributed Problems. *Int. Journal of Comp. and Inf. Sci.* 13. 1984, Pages 131-147.