

TOPOLOGY ABSTRACTION AND OVERLAYS

- Physical topology
nodes represent all network nodes (including switching elements, non-app executing ^{hosts} nodes)
edges represent all communication links
- Topology overlay
• Superimposed topology
higher level topology superimposed on logical topology. usually a regular structure like tree, ring, mesh, or hypercube. it provides a specialized path for efficient information dissemination and/or gathering as part of a distributed algorithm.
- application execution
execution of distributed application program.
- control algorithm (protocol)
executed in order to monitor application execution, or perform various auxiliary functions: creating a spanning tree, creating a connected dominating set, achieving consensus among nodes, dist. transaction commit, dist. deadlock detection, global predicate detection, termination detection, ...

- logical topology
nodes represent all end hosts where application executes.
edges are logical channels among these nodes

- Centralized algorithm
most work is performed by one (or few) processors, whereas others play a relatively small role in performing the joint task.
- distributed algorithm
each processor plays an equal role in sharing the msg, time and space overhead.
- symmetric algorithm
all processors execute the same logical functions.
- asymmetric algorithm
different processors execute logically different functions.

{Structural elegance}

- anonymous algorithm

does not use process identifiers in code.

{allows scalability + transparency}

- Uniform algorithm

does not use n (total no. of processes) as parameter in code.

- adaptive algorithm

complexity of the algorithm can be expressed in terms of k (no. of participating nodes).

- deterministic receive primitive

specifies the source from which it wants to receive a message.

- non-deterministic receive primitive

can receive a msg from any source.

- deterministic execution

program that contains no non-deterministic receives.

- non-deterministic execution

program that contains at least one non-deterministic primitive.

- execution inhibition

protocols that require processors to suspend their normal execution until some action has been performed are called inhibitory or freezing protocols.

- locally, globally inhibitory

- send, receive, internal event inhibition

- Synchronous system

system that has known upper bound in:

- message delay, - clock drift rate,

- logical step execution time.

- asynchronous system

system which satisfies none of above 3 conditions.

- on-line algorithm

executes as data is being generated.

- off-line algorithm

requires all data to be available before execution begins.

- wait-free algorithm

can execute in an $(n-1)$ process fault tolerant manner.

(resistant to $n-1$ process failures)

PROCESS FAILURE MODELS

- fail-stop

Execution stops at some instant. other processes learn of this failure

- crash

Execution stops at some instant. other processes do not learn of this failure.

- receive omission

intermittently fail to receive some messages.

- send omission

intermittently fail to send some messages.

- general omission

intermittently fail to receive / send some msgs.

- byzantine or malicious failure, with authentication

processes may exhibit arbitrary behaviour, but signatures can be used to verify msg integrity.

- byzantine or malicious failure

processes may exhibit arbitrary behaviour.
(no authentication techniques applicable)

COMMUNICATION FAILURE MODELS

- crash

link stops carrying msgs at some instant.

- omission

intermittently fails to carry some messages.

- byzantine failure

link can exhibit arbitrary behaviour, including creating spurious messages, and modifying the msgs sent on it.

- timing failures

these can occur in synchronous systems either for processes, or for links.

- complexity measures & metrics

- space, time & node vs. systemwide

- message complexity

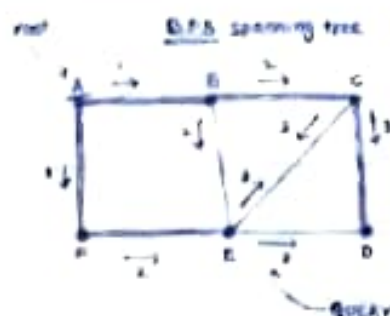
number, size, time.

SYNCHRONOUS SINGLE-INITIATOR SPANNING TREE ALGORITHM USING FLOODING

each node already knows its neighbours. spanning tree is known, when each node knows its parent. this algorithm assumes a designated root node.

root: • starts by sending QUERY to its neighbours.

nodes: • on receiving QUERY(s) one random node is chosen as parent. QUERY is sent to all remaining neighbours.



6 nodes & 8 links

3 rounds 9 QUERYs

Calgo uses "visited", and "depth" variables.
(diameter = longest distance in graph)

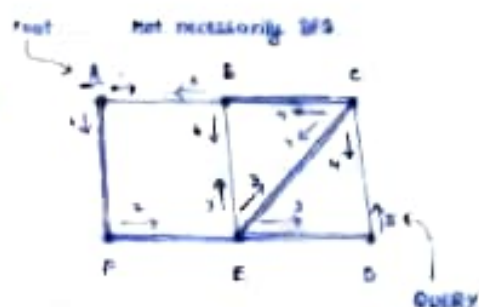
ASYNCHRONOUS SINGLE-INITIATOR SPANNING TREE ALGORITHM USING FLOODING

this algorithm assumes a designated root node. algorithm terminates when response to all QUERYs is received.

roots: • starts by sending QUERY to its neighbours.

nodes: • on receiving first QUERY, it is chosen as parent.

and ACCEPT reply is sent to it. to remaining QUERY mgs, REJECT reply is sent. QUERY is sent to all remaining neighbours.



ASYNCHRONOUS CONCURRENT-INITIATOR SPANNING TREE ALGORITHM USING FLOODING

suppose the any node may spontaneously initiate spanning tree algorithm provided it has not already been invoked locally due to receipt of a QUERY message

Suppress the instance initiated by one root and continue the instance initiated by the other root, based on tie-breaking using the processor identifier. this algorithm uses $QUERY(\text{newroot})$, $ACCEPT(\text{newroot})$, and $REJECT(\text{newroot})$ messages.

Root: starts by sending $QUERY(i)$ to its neighbours.

node: when $QUERY(\text{newroot})$ received from j :

↳ $\text{newroot} > \text{myroot}$: suppress current execution. reset data structures.

join j 's subtree with $\text{myroot} = \text{newroot}$.

↳ $\text{newroot} = \text{myroot}$: already have a parent from same root.

hence send $REJECT$ to j .

↳ $\text{newroot} < \text{myroot}$: j 's root has lower priority, so send $REJECT$ to j .

j will eventually receive $QUERY(\text{myroot})$ from i .

when $ACCEPT(\text{newroot})$ received from j :

↳ $\text{newroot} = \text{myroot}$: $ACCEPT$ is in response to $QUERY$ sent by i .

process normally.

↳ $\text{newroot} < \text{myroot}$: $ACCEPT$ is response to old $QUERY$ from i (with old root).

ignore message.

↳ $\text{newroot} > \text{myroot}$: cannot happen.

when $REJECT(\text{newroot})$ received from j : similar to $ACCEPT$ above.

ACCEPTED REJECTED
when children U unrelated = neighbours - parent; send $ACCEPT(\text{myroot})$
to parent

drawback: only root knows when the algorithm has terminated. It must broadcast others to inform this. (along spanning tree).

ASYNCHRONOUS CONCURRENT-INITIATOR DEPTH FIRST SEARCH SPANNING TREE ALGORITHM

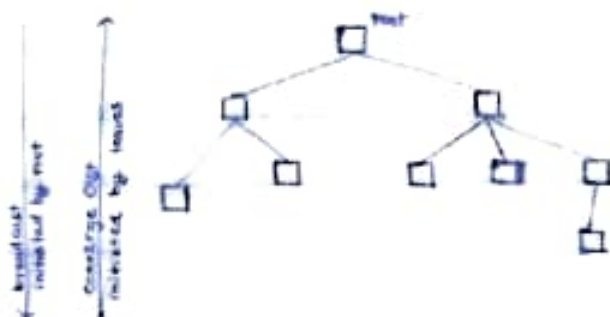
similar to the algorithm before, though here QUERY is sent to only one neighbour at a time after the receipt of a QUERY, ACCEPT, or REJECT.

BROADCAST AND CONVERGECAST ON A TREE

spanning tree is useful for distributing info to all nodes via - broadcast

and collecting info from all nodes via

- convergecast ← leader election



SINGLE SOURCE SHORTEST PATH ALGORITHM: SYNCHRONOUS BELLMAN-FORD



after k rounds, length of nodes k hops away from source is stabilized.

longest path length = $n-1$

rounds to stabilize = $n-1$

initially, length = ∞

for source, length = 0

each node sends its length to all its neighbours.

on receiving length, the node updates its length if

$\text{length}_i + \text{weight}_i < \text{length}_i$

initially ∞

weight between nodes and receiver

msg complexity = $O((n-1)E)$

links

algorithm assumes no cyclic paths having negative weight.

SINGLE-SOURCE SHORTEST PATH ALGORITHM: ASYNCHRONOUS BELLMAN-FORD

Very similar to previous algorithm where, a node terminates after sending UPDATE to all neighbours. Unfortunately msg complexity: $\Theta(n^2)$, time comp: $\Omega(n^2 \cdot d)$

MINIMUM-WEIGHT SPANNING TREE (MST) ALGORITHM IN A SYNCHRONOUS SYSTEM

Kruskal's algorithm: begins with a forest of graph components. in each iteration it identifies the min weight edge that connects 2 different components, and uses this edge to merge 2 components. this continues until all components are merged into single component.

prim's algorithm & dijkstra's algorithm: a single node component is selected. in each iteration, a min weight edge incident on the component is identified, and the component expands to include that edge and the node at the other end of that edge. after $n-1$ iterations, all the nodes are included.

distributed algorithm by gallagher, humblet, and spira (GHS) uses kruskal's strategy.

MWOE - min weight outgoing edge

this algorithm knows the value of n . it assumes that weight of each edge in the network is unique, which is necessary to guarantee unique MST. (if weights are not unique, the ids of the nodes on which they are incident can be used as tie-breakers by defining a well-formed order)

localid $\xrightarrow{\text{MWOE}}$ remoteid

- initially each component has one node which is the leader.
 - leader broadcasts $\text{SEARCH_MWOE}(\text{leader})$ on tree edges to request edges to find x .
 - edges send $\text{EXAMINE}(\text{leader})$ to non-tree edges after receiving SEARCH_MWOE .
 - details of potential MWOEs are convergecast to leader with $\text{REPLY_MWOE}(\text{localid}, \text{remoteid})$.
 - leader broadcasts $\text{ADD_MWOE}(\text{localid}, \text{remoteid})$ to add MWOE and identify new leader.
 - new leader is $\max(\text{localid}, \text{remoteid})$ which touches the MWOE.
 - new leader broadcasts $\text{NEW_LEADER}(\text{leader})$ after merging components.
 - a cycle of length 2 is possible, but since they have a common MWOE, there would be a single spanning tree, and a unique leader in merged comp.
 - a cycle of length 2 or more is not possible as that could mean 2 different minimum weight edges to one tree, which is not possible.
- iteration time complexity = $O(\log n)$ total time complexity = $O(n \log n)$
 n msgs to tree edges, 2 msgs to find MWOE. \rightarrow msg complexity = $O(n \log n)$

MAXIMAL INDEPENDENT SET (MIS)

proposed by Luby

each iteration:

- each process selects a random number, random_i .
 - if its random_i is least among neighbours it gets to be in MIS.
 - each P_i informs its selection status through $\text{SELECTED}(\text{yes})$ msg.
 - if a process finds atleast one of its neighbours is $\text{SELECTED}(\text{yes})$ it eliminates itself from MIS.
 - each P_i informs its elimination status through $\text{ELIMINATED}(\text{yes})$ msg.
 - each P_i removes from its neighbours, who send $\text{ELIMINATED}(\text{yes})$ msg.
 - proceed to next iteration.
- no of iterations = $O(\log n)$

