#### Distributed Algorithms for Computer Networks

## Chapter 5 Spanning Tree Construction

Dr.Ahmed Awad Dr. Amjad Hawash

#### Introduction

- **Spanning Tree**: The maximal set of edges of a graph that contains no cycles.
- → The minimal set of edges that connect all vertices.
- Why spanning trees in computer networks?
  - They provide a subtree with possibly less communication links, resulting in lowered communication costs.
- Providing a parent/child relationship among the nodes of the network eases the task of communication since the source and the destination of the communication is known beforehand.

### The Flooding Algorithm

- **Broadcast:** Sending a message to all nodes in a computer network.
- **Flooding**: Forward any incoming message to all neighbors except the neighbor that has sent the message.
- If the same message arrives again, it should be discarded.
- The node that initiates the message to be broadcasted is called the **root**.

### The Flooding Algorithm (2)

Alg	Algorithm 4.1 Flood			
1:	<b>int</b> <i>i</i> , <i>j</i>			
2:	<b>boolean</b> visited $\leftarrow$ false			
3:	message types msg			
4:				
5:	if $i = root$ then	▷ <i>root</i> initiates flooding		
6:	send <i>msg</i> to $\Gamma(i)$			
7:	$visited \leftarrow true$			
8:	end if			
9:				
10:	<b>receive</b> <i>flood</i> ( <i>j</i> )	⊳ <i>flood</i> may be received many times		
11:	<b>if</b> <i>visited</i> $=$ <i>false</i> <b>then</b>	$\triangleright$ msg received first time		
12:	send <i>msg</i> to $\Gamma(i) \setminus \{j\}$			
13:	$visited \leftarrow true$			
14:	else	▷ <i>msg</i> received before		
15:	discard msg			
16:	end if			

The Flooding Algorithm : Complexity Analysis

#### • Message Complexity:

- Each edge connects two nodes.
- Each edge is used to deliver a message at least once and at most twice when two nodes send messages concurrently.
- Total number of messages =2m where m represents the number of edges in the graph.
- → Message complexity: O(2m)=O(m)

#### • Time Complexity:

- The longest time for the broadcast message to reach any node in the graph is the distance between the two farthest nodes of the graph.
- $\rightarrow$  Time complexity: *d*: The diameter of the graph.

#### Flooding-Based Asynchronous Spanning Tree Construction Algorithm (Flood\_ST)

- **Objective:** Build spanning tree originating from the initiator root for broadcasting. This tree can be used for any further broadcast messages.
- Assumption: Each node in the tree except the leaf nodes should know the identifiers of its children and all nodes except the root should know their parents.
- Any node that wants to build a broadcast tree initiates the algorithm and becomes the root of the spanning tree to be formed.

### Flood\_ST

#### • Types of messages:

- Probe.
- Ack.
- Reject
- Any node that wants to build a spanning tree starts the algorithm by sending **probe** message to its neighbors.
- When a node receives a probe message:
  - it sends a positive acknowledgement to the source node (which becomes a parent) in case it doesn't have a parent.
  - It transfers the probe message to all of its neighbors except the parent.
  - If a node has a parent, it sends a reject message to source node of the probe message.

### Flood\_ST Algorithm

#### Algorithm 4.2 Flood\_ST

```
1: int parent \leftarrow \perp
 2: set of int childs \leftarrow \emptyset, others \leftarrow \emptyset
 3: message types probe, ack, reject
 4:
 5: if i = root then
                                                                              \triangleright root initiates tree construction
 6:
          send probe to \Gamma(i)
 7:
          parent \leftarrow i
 8: end if
 9:
10: while (childs \cup others) \neq (\Gamma(i) \setminus \{parent\}) do
11:
          receive msg(j)
12:
          case msg(j).type of
                                                                                     ▷ probe received first time
13:
                       probe:
                                   if parent = \perp then
14:
                                       parent \leftarrow j
                                       send ack to j
15:
16:
                                       send probe to \Gamma(i) \setminus \{j\}
17:
                                    else
                                                                                         ⊳ probe received before
18:
                                       send reject to j
19:
                                    childs \leftarrow childs \cup \{j\}
                                                                                          \triangleright include j in children
                       ack:
20:
                                   others \leftarrow others \cup \{j\}
                                                                          \triangleright include j in unrelated neighbors
                       reject:
21: end while
```

### Flood\_ST- Example



Message complexity= 2m+2=20 Time complexity=3

→ Shortest path is not always followed due to communication link speeds.

### Flood\_ST: Complexity Analysis

#### Message Complexity:

- Each edge will be traversed at least twice with probe and ack or probe and reject messages.
- Each edge will be traversed at most 4 times in the case of two nodes attempting to send each other probe message concurrently.
- In total 4m messages will be sent.
- → Message complexity = **O(4m)=O(m)**

### Flood\_ST: Complexity Analysis

#### • Time Complexity:

- The message maybe transferred over the fast communication links of the longest path instead of the shortest paths with slow links.
- → Time Complexity = O(n)
   n: The number of vertices in the graph.



#### Breadth First Search (BFS)

- **Graph Traversal**: Visiting all vertices of a graph in some predefined order.
- It maybe required to find the shortest distances of vertices from a source vertex in terms of the number of links (hops) to that source.
- → Approach: Breadth-First-Search (BFS) tree.

**Definition 5.1** (Breadth-First-Search Tree) A *breadth-first-search tree* T of a graph G is a spanning tree of G such that for every node of G, the tree path is a minimum-hop path to the root.

#### Breadth First Search (BFS)





G

E

Enqueued Array А B  $\mathbf{Q} \rightarrow$ С D E F G Η

How is this accomplished? Simply replace the stack with a queue! Rules: (1) Maintain an *enqueued* array. (2) Visit node when *dequeued*.



Enqueued Array Α В  $Q \rightarrow D$ С  $\sqrt{}$ D E F G Η

**Nodes visited:** 

**Enqueue D. Notice, D not yet visited.** 





Enqueued Array А В  $Q \rightarrow C \rightarrow E \rightarrow F$  $\sqrt{}$ С D  $\sqrt{}$ Е  $\sqrt{}$ F  $\sqrt{}$ G Η

**Nodes visited: D** 

Dequeue D. Visit D. Enqueue unenqueued nodes adjacent to D.



#### Dequeue C. Visit C. Enqueue unenqueued nodes adjacent to C.



Enqueued  $Q \rightarrow F \rightarrow G$  $\sqrt{}$  $\sqrt{}$  $\sqrt{}$  $\sqrt{}$ Η

**Nodes visited: D, C, E** 

Dequeue E. Visit E. Enqueue unenqueued nodes adjacent to E.





Enqueued Array А В  $Q \rightarrow G$  $\sqrt{}$ С D  $\sqrt{}$ Ε  $\sqrt{}$ F  $\sqrt{}$ G  $\sqrt{}$ Η

Nodes visited: D, C, E, F

Dequeue F. Visit F. Enqueue unenqueued nodes adjacent to F.



Dequeue G. Visit G. Enqueue unenqueued nodes adjacent to G.



#### Dequeue H. Visit H. Enqueue unenqueued nodes adjacent to H.



Α  $\sqrt{}$ В  $\sqrt{}$  $Q \rightarrow B$  $\sqrt{}$ С D  $\sqrt{}$ Е  $\sqrt{}$ F  $\sqrt{}$ G  $\sqrt{}$ Η  $\sqrt{}$ 

Nodes visited: D, C, E, F, G, H,

A

**Dequeue A. Visit A. Enqueue unenqueued nodes** adjacent to A.







Nodes visited: D, C, E, F, G, H, A, B

Dequeue B. Visit B. Enqueue unenqueued nodes adjacent to B.



**Q** empty. Algorithm done.

### Depth First Search (DFS)

- **Objective**: Find all of the vertices reachable from a source vertex in the graph.
- It visits all possible vertices as far as it can reach and when all vertices are visited, it returns to the parent node.

**Definition 5.2** (Depth-First-Search Tree) A *frond* edge is an edge that does not belong to a spanning tree. For a rooted spanning tree T of a graph G, let us denote by S(u) all the nodes in the subtree of u, and P(u) denote all the vertices that exist between u and the root. A *depth-first-search tree* of a graph G is a spanning tree T of G such that for every frond edge  $\{u, v\}, v \in S(u) \lor v \in P(u)$  [5].

- Applications:
  - Finding strongly connected components of a directed graph.
  - Cycles detection !!

.

#### Depth First Search (DFS)





## Task: Conduct a depth-first search of the graph starting with node D



Visited Array			
	A		
	В		
	С		
	D	$\checkmark$	
	E		
	F		

G

Η

The order nodes are visited:

Visit D

D

D



The order nodes are visited:

D



Consider nodes adjacent to D, decide to visit C first (Rule: visit adjacent nodes in alphabetical order)



Visited				
Α	Array			
	Α			
	В			
	С	$\checkmark$		
	D	$\checkmark$		
	E			
	F			
	G			
	Η			



Visit C

The order nodes are visited:

D, C



The order nodes are visited:

D, C

Visited Array



No nodes adjacent to C; cannot continue → backtrack, i.e., pop stack and restore previous state



The order nodes are visited:

D, C

Visited Array



#### Back to D – C has been visited, decide to visit E next



The order nodes are visited:

D, C, E

Visited Array



#### Back to D – C has been visited, decide to visit E next



Visited Array Α В  $\sqrt{}$ С  $\sqrt{}$ D  $\sqrt{}$ E F E G Η D

The order nodes are visited:

D, C, E

Only G is adjacent to E



Visited Array



The order nodes are visited:

Visit G

D, C, E, G



The order nodes are visited:

D, C, E, G

Visited Array



Nodes D and H are adjacent to G. D has already been visited. Decide to visit H.



Visited Array



The order nodes are visited:

Visit H

D, C, E, G, H



The order nodes are visited:

D, C, E, G, H

Visited Array



Nodes A and B are adjacent to F. Decide to visit A next.



Visited Array  $\checkmark$ Α В  $\sqrt{}$ С  $\sqrt{}$ D  $\sqrt{}$ E F  $\sqrt{}$ G Η  $\sqrt{}$ 

Visit A

А

Η

G

E

D

The order nodes are visited:

D, C, E, G, H, A



Visited Array



The order nodes are visited:

D, C, E, G, H, A

Only Node B is adjacent to A. Decide to visit B next.



Vi	Visited		
A	Array		
	А	$\checkmark$	
	В	$\checkmark$	
	С	$\checkmark$	_
	D	$\checkmark$	_
	Е	$\checkmark$	-
	F		_
	G	$\checkmark$	
	Η	$\checkmark$	
		1	

Visit B

В

А

Η

G

E

D

The order nodes are visited:

D, C, E, G, H, A, B



The order nodes are visited: D, C, E, G, H, A, B Visited Array



#### No unvisited nodes adjacent to B. Backtrack (pop the stack).



The order nodes are visited: D, C, E, G, H, A, B Visited Array



No unvisited nodes adjacent to A. Backtrack (pop the stack).



The order nodes are visited: D, C, E, G, H, A, B Visited Array



No unvisited nodes adjacent to H. Backtrack (pop the stack).



Visited Array



The order nodes are visited:

D, C, E, G, H, A, B

No unvisited nodes adjacent to G. Backtrack (pop the stack).



The order nodes are visited: D, C, E, G, H, A, B



# No unvisited nodes adjacent to E. Backtrack (pop the stack).



Visited Array



The order nodes are visited:

D, C, E, G, H, A, B

F is unvisited and is adjacent to D. Decide to visit F next.



Visited Array  $\checkmark$ Α

В

С

D

E

F

G

Η

 $\sqrt{}$ 

 $\sqrt{}$ 

 $\sqrt{}$ 

 $\sqrt{}$ 

 $\sqrt{}$ 

 $\sqrt{}$ 



The order nodes are visited:

D, C, E, G, H, A, B, F







The order nodes are visited:

D, C, E, G, H, A, B, F

No unvisited nodes adjacent to F. Backtrack.





The order nodes are visited: D, C, E, G, H, A, B, F

No unvisited nodes adjacent to D. Backtrack.



Visited Array



The order nodes are visited:

D, C, E, G, H, A, B, F

Stack is empty. Depth-first traversal is done.

#### Distributed Breadth-First-Search Algorithms

- **Synchronous BFS**: Synchronous algorithm working in rounds to construct the BFS tree.
- Asynchronous BFS Construction: Asynchronous algorithm to construct the BFS tree.

#### Synchronous BFS Construction Algorithm (Synch\_BFS)

- The synchronous distributed version of Dijkstra's algorithm for single-source shortest path problem.
- A single initiator algorithm.
- In each synchronous round, a partial BFS tree is formed.
- The already formed branches of the tree are used to carry synchronization messages and the leaves search for new nodes to be added to the tree.
- The depth of the tree is incremented in each round until all the nodes are covered.

### Synch\_BFS Algorithm Features

- Full termination detection is provided so that all other nodes know when the BFS is constructed.
- Synchronization is performed using special control messages
  - ➔ Eliminates the need of other synchronization techniques such as specific synchronizer node.

### Synch\_BFS Algorithm Message Types

- Round: Sent by the root at the beginning of each round and transferred by all the nodes of the partial BFS tree to their children.
- **Probe:** Sent by the leaves of the partial BFS tree to the unsearched neighbors (that have not yet been members of the tree).
- Ack/Reject: Sent by the searched node to accept/reject being a child of the sending leaf node.
- **Upcast:** Sent first by the leaf nodes of the partial BFS tree and then by the intermediate nodes to their parents to signal the end of neighbor search and the end of the round.
- Finish: Sent by the leaf nodes to their parents to signal that their part is done as they either have no neighbors other than the parent or they do not have any children.
- **Terminate:** Broadcasted by the root to all nodes to signal that the construction of the BFS tree is

### Synch\_BFS Algorithm Flow

→The root starts the algorithm by sending the first round message to its neighbors.

→The neighbors respond with **ack** message.

➔ The round message is transferred over the partial BFS tree to the leaf nodes.

→ The leaf nodes search nodes for the next level the BFS tree by sending probe messages.

→ Each leaf node that has received ack or reject message from all of its neighbors except the parent returns upcast message to its parent, which convergecasts the upcast message to its parent.

→ When all upcast messages from the neighbors of the root are received, root starts the next round by issuing the next round message.

#### Synch\_BFS Algorithm Flow- Example



### Synch\_BFS Algorithm Termination

- The root does not know beforehand the diameter of the graph.
- When all the neighbors of a node i except its parent have responded by a reject message or when a leaf node ( in the constructed BFS tree) does not have any neighbors other than its parent, the part of node i is over.
- When a node reaches this situation, it sends a finish message to its parent which converecasts it to the root.
- When the root receives finish message from all of its children, it broadcasts a **terminate** message over the network to inform every node that the process is over.

Algorithm 5.1 Synch\_BFS

1:	<b>int</b> <i>parent</i> $\leftarrow \perp$ , $k \leftarrow 1$ , $i, j$			
2:	set of int childs, others, phased, finished $\leftarrow \varnothing$			
3:	message types round, probe, ack, reject, upcast			
4:	<b>boolean</b> $leaf_flag \leftarrow false$			
5:				
6:	if $i = root$ then			
7:	send $probe(k)$ to $\Gamma(i)$			
8:	while true do			
9:	<b>receive</b> $msg(j)$			
10:	case msg.type of			
11:	$ack(k) \lor upcast(k)$ :	$phased \leftarrow phased \cup \{j\}$		
12:	finish:	finished $\leftarrow$ finished $\cup \{j\}$		
13:		if finished = childs then		
14:		send terminate to ch	ilds	
15:		exit	⊳ root terminates	
16:	if phased = childs then		⊳ start next round	
17:	$k \leftarrow k+1$			
18:	<b>send</b> $round(k)$ to $ch$	hilds, phased $\leftarrow 0$		
19:	end while			

20:	else		
21:	while true do		
22:	$round\_over \leftarrow false$		
23:	<b>while</b> ¬ <i>round_over</i> <b>do</b>		
24:	<b>receive</b> $msg(j)$		
25:	case msg.type of		
26:	round(k):	if $state = interm$ then	▷ intermediate node
27:		<b>send</b> $round(k)$ to <i>childs</i>	
28:		else	⊳ leaf node
29:		<b>send</b> $probe(k)$ to $\Gamma(i) \setminus \{parentifyed)$	ent}
30:		$round\_recvd \leftarrow true$	
31:	probe(k):	if $parent = \perp$ then	⊳ non-member node
32:		parent $\leftarrow j$ ; state $\leftarrow$ new_lea	af
33:		<b>send</b> $ack(k)$ to $j$	
34:		else	⊳ a member node
35:		<b>send</b> $reject(k)$ to $j$	
36:	ack(k):	$childs \leftarrow childs \cup \{j\}$	
37:		if state $\neq$ interm then interm_flag	$r \leftarrow true$
38:	reject(k):	others $\leftarrow$ others $\cup \{j\}$	
39:		if others = $\Gamma(i) \setminus \{parent\}$ then	$\triangleright$ a leaf node finishes
40:		send finish to parent	
41:	upcast(k):	$phased \leftarrow phased \cup \{j\}$	
42:	finish:	finished $\leftarrow$ finished $\cup \{j\}$	
43:		<b>if</b> ( <i>finished</i> = <i>childs</i> ) $\land$ ( <i>state</i> = <i>in</i> )	tterm) then
44:		send finish to parent	
45:	<u>terminate</u> :	if $childs \neq 0$ then	
46:		send terminate to childs	
47:		exit	⊳ all nodes terminate

48:	if round_recvd then		
49:	if $((state = leaf) \land (($	childs $\cup$ othe	$rs) = (\Gamma(i) \setminus \{parent\})) \lor ((state =$
	$interm$ ) $\land$ (childs = phased $\cup$ finished	)) then	$\triangleright$ check end of round
50:	<b>send</b> <i>upcast(p)</i> to <i>p</i>	arent	
51:	<b>if</b> $state = new\_leaf$	then	⊳ update states
52:	$state \leftarrow leaf$		
53:	else if interm_flag t	hen	
54:	$state \leftarrow interm$		
55:	end if		
56:	phased, others $\leftarrow \&$	; round_rec	$vd \leftarrow false; round\_over \leftarrow true$
57:	end if		
58:	end if		
59:	end while		
60:	end while		
61:	end if		

Synch\_BFS Algorithm Complexity Analysis

- **Lemma**: Synch\_BFS algorithm correctly constructs a BFS tree.
  - **Proof**: By **mathematical induction**, let k denote the step at which the BFS tree is being constructed.
    - (1) At step k=0, no tree is constructed.
    - (2) Assume that at step k,  $T_k$  tree rooted at r is constructed.
    - (3) at step k+1, only the leaves of the tree will be active in sending the probe messages to their neighbors that are one hop away.
- Any added nodes will be k+1 hops away from r, forming T`<sub>k</sub>

Synch BFS Algorithm Complexity Analysis

- Time Complexity:
   Broadcasting the round message to current T<sub>k</sub> BFS tree (whose depth is k) is performed in k units.
  - Convergecast of the upcast messages is performed in k units as well.
  - Additional time is required for probe and ack/reject messages.
  - $\rightarrow$  For each step, 2k+2 time is required.
  - The depth of the tree will be d in total (d represents the diameter of the graph).

• Time Complexity = 
$$2\sum_{k=1}^{d} (k+1) = (d+1)(d+2)/2 = O(d^2).$$

#### $\rightarrow O(d^2)$

d: The diameter of the graph.

Synch\_BFS Algorithm Complexity Analysis

#### Message Complexity:

- In round p, if the tree formed has k nodes, there will be k-1 round messages and k-1 upcast messages.
- Total number of synchronization messages for a round is O(n), n: the number of vertices in the graph.
- The number of synchronization messages in all rounds will be O(nd)
- In each round, new edges of the BFS tree will be determined by probe and ack/reject messages.
- The total number of discovery messages=2m

m: The number of edge O(nd) + O(m) = O(nd + m)

➔ Message Complexity=