

Distributed Algorithms ~~& book by Donald~~

- Global state, global time, non-determinism.
- Application: Telecommunication, distributed information processing, scientific computing, real time process control.
- Topology, undirected graph.
Three levels of typology: physical, logical, superimposed.

Classification and basic concepts

- Application execution and control algorithm execution (Protocol).
- centralized algorithm and distributed algorithms.

23/10/12.

- Semester Symmetric and Asymmetric algorithms.
- Anonymous algorithms: An anonymous system is a system in which neither processes nor processors use their identifiers to make any execution decisions in the distributed algorithm. An anonymous algorithm is an algorithm which runs on an anonymous system in which and therefore does not use process or processor identifiers in the code. In these we have structured elegance and these are hard to design.
- Uniform algorithms: A uniform algo. does not use n , the no. of processes in the system as a parameter in its code. These allow scalability transparency.
- Adaptive algorithm: Consider the context of a problem X .
In a system with n nodes, let k , $k \leq n$ be the no. of nodes participating in the context of X . If complexity of the system can be expressed in terms of k rather than n , the algorithm is adaptive.

- Deterministic v/s Non-deterministic execution
- Execution inhibition : blocking communicating primitives freeze the local execution. Non blocking, inhibition or freezing protocols. Global inhibiting - send/receive inhibitory, internal event inhibitory.
- Synchronous / Asynchronous Systems :
For synch. ~~as~~ systems -
 - i) known upper bound on the msg communication delay
 - ii) known bounded drift rate for the clock.
 - iii) known upper bound to execute a logical step in execution.
- Online v/s offline algorithms : based on data.
- Failure Models : t-fault tolerant
 - Process failure models : fail-stop, crash, receive omission, send omission, general omission,
 - Byzantine failure models or malicious failures with auth, Byzantine or malicious failures
 - Communication failure models : crash failure, omission failures, byzantine failures
- wait-free algorithms : resilient, fault-tolerant
- communication channels : FIFO or non-FIFO.
- complexity measures and metrics :
For sequential algorithms : time & space complexity, lower bound (Ω, ω), upper bound (O, o)

and exact bound (Θ).

For distributed algorithms : time & space complexity, message complexity, space-complexity per node, system-wide space complexity, time-complexity per node, system-wide time complexity, msg complexity (no. of msgs, and size of msgs and msg time complexity).

- Program Structure : Hoare, CSP (concurrent sequential processes)
Proposed by

$$* [G_1 \rightarrow C_{L1} \parallel G_2 \rightarrow C_{L2} \parallel \dots G_{Lk} \rightarrow C_{Lk}]$$

26/10/19.

Synchronous single-initiator spanning tree algorithm using n flooding
(Kshem)

Synchronous BFS spanning tree algorithm.

The code shown is for process P_i , $1 \leq i \leq n$.
(local variables).

int visited, depth $\leftarrow 0$.

int parent $\leftarrow 1$

set of int Neighbors \leftarrow set of neighbors.

(message type)

QUERY.

if $i = \text{root}$ then

visited $\leftarrow 1$,

depth $\leftarrow 0$

send QUERY to Neighbors;

for round = 1 to diameter do

if visited = 0 then

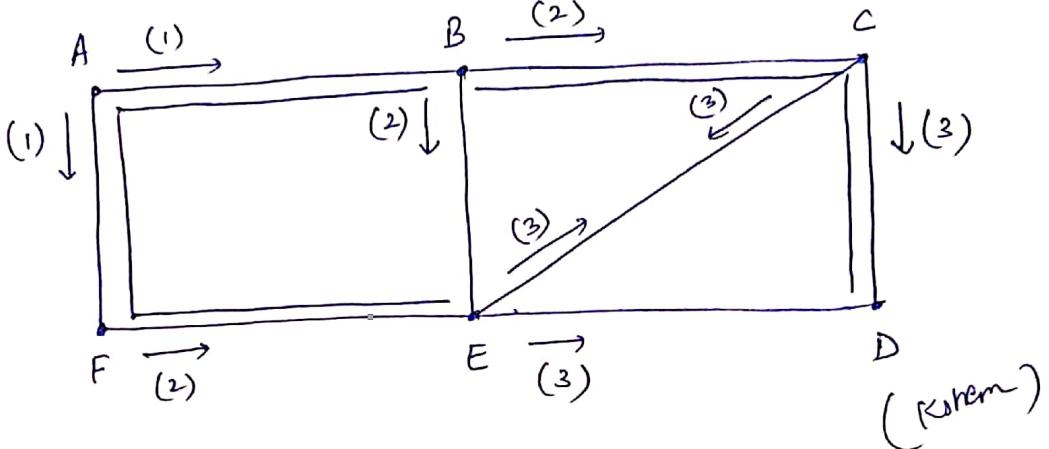
if any QUERY messages arrive then

parent \leftarrow randomly select a node from which QUERY was received.

$\text{visited} \leftarrow 1$
 $\text{depth} \leftarrow \text{round}$
Send QUERY to Neighbors \ {senders of QUERYS
 received in this round}.

delete any QUERY msgs that arrived in this round.

Example.



Asynchronous single-initiator spanning tree algorithm

The code for process P_i , $1 \leq i \leq n$

(local variables)

int parent $\leftarrow \perp$ \perp \leftarrow undefined

set of int children, unvisited $\leftarrow \emptyset$

set of int neighbors \leftarrow set of neighbors

(Message types)

QUERY, ACCEPT, REJECT

When the predesignated root node wants to initiate the algo.

if ($i = \text{root}$ and $\text{parent} = \perp$) then

Send QUERY to all neighbors

parent \leftarrow i

when QUERY arrives from j :

```

if parent = 1 then
    parent = j
    send ACCEPT to j;
    send QUERY to all neighbors except j;
    if (children ∪ unrelated) = (Neighbors \ {parent})
        then terminate;

```

when ACCEPT arrives from j:

`children ← children ∪ {j};`

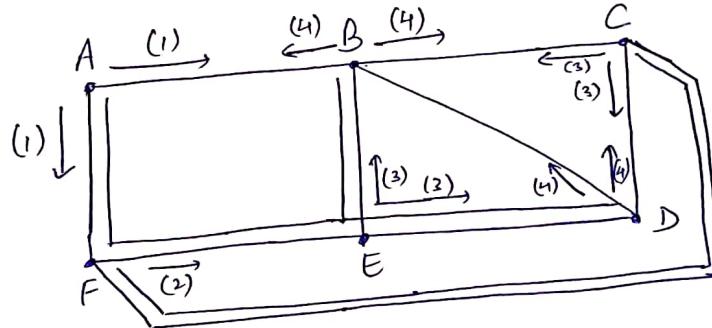
$\text{children} \leftarrow \text{children} \cup j$,
if $(\text{children} \cup \text{unvisited}) = (\text{Neighbors} \setminus \{\text{parent}\})$ then terminate;

When REJECT arrives from j :

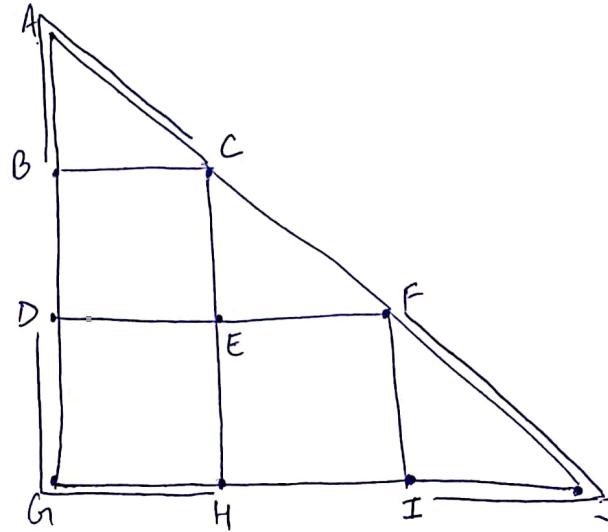
unrelated \leftarrow unrelated $\cup \{j\}$;

$\text{unrelated} \leftarrow \text{unrelated} \cup \{j\}$
 if $(\text{children} \cup \text{unrelated}) = (\text{Neighbours} \setminus \{\text{parent}\})$ then terminate;

Example.



Asynchronous concurrent-initiator spanning tree algorithm using flooding
An (Kishem)



(local variables)

int parent, myroot $\leftarrow 1$

set of int children, unrooted $\leftarrow \emptyset$

set of neighbors \leftarrow set of neighbors

(message types)

QUERY, ACCEPT, REJECT

When the node wants to initiate the algorithm as a root :

if (parent = 1) then

send QUERY(i) to all neighbors;

 parent, myroot $\leftarrow i$

when QUERY(newroot) arrives from j :

if (myroot < newroot) then ~~discard earlier partial execution~~

 // discard earlier partial execution due to its low priority.

then parent $\leftarrow j$;

 myroot \leftarrow newroot;

 unrooted $\leftarrow \emptyset$;

send QUERY(newroot) to all neighbors except j;

if Neighbors = {j} then // leaf node .

send ACCEPT(myroot) to j;

 terminate;

else send REJECT(newroot) to j;

 // if newroot = myroot then parent is already identified

 // if newroot < myroot then ignore the QUERY, j will update

 // its root when it receives QUERY(myroot) .

When ACCEPT(newroot) arrives from j :

if (newroot = myroot) then

children \leftarrow children $\cup \{j\}$

if (children \cup unrelated) = (Neighbors \ {parent}) then ~~terminate~~

if (i = myroot) then terminate

else send ACCEPT(myroot) to parent.

// if newroot < myroot then ignore the msg and newroot > myroot

// will never occur.

when REJECT(newroot) arrives from j :

if (newroot = myroot) then

unrelated \leftarrow unrelated $\cup \{j\}$.

if (children \cup unrelated) = (Neighbors \ {parent}) then

if (i = myroot) then terminate

else send ACCEPT(myroot) to parent.

// if newroot < myroot then ignore the msg and newroot >

// myroot will never occur.

30/10/19.

Asynchronous concurrent-initiator DFS spanning tree algorithm

(local variables)

int parent, myroot \leftarrow 1;

set of int children, unrelated $\leftarrow \emptyset$;

set of int Neighbors, unknown \leftarrow set of neighbors

(message types)

QUERY, ACCEPT, REJECT

when the node wants to initiate the algorithm as root.

if (parent = 1) then send QUERY(i) to i (itself).

when QUERY(newroot) arrives from j :

if myroot < newroot then

parent \leftarrow j; myroot \leftarrow newroot, unknown \leftarrow set of neighbors;

unknown \leftarrow unknown \ {j};

if unknown \neq 0 then

delete some x from unknown.

send QUERY(myroot) to x.

else send ACCEPT(myroot) to j;

else if (myroot = newroot) then

~~send~~ REJECT(myroot) to j;

// if newroot < myroot ignore the query.

// j will update its root to higher root identifier when it receives

// its QUERY.

when ACCEPT(newroot) or REJECT(newroot) arrives from j :

if (newroot = myroot) then

if ACCEPT arrived then

children \leftarrow children \cup {j}.

if (unknown = 0) then

if parent \neq i then

send ACCEPT(myroot) to parent;

else set i as the root; terminate

else delete some x from unknown;

send QUERY(myroot) x;

// if newroot < myroot ignore the query. Since sending QUERY to
// j, i has updated its myroot. j will update its myroot to a
// higher root identifier. When it a query initiated by it
// newroot > myroot will never occur.

* Michael T. Goodrich and Roberto Tamassia "Algorithm design".

The Bellman-Ford shortest path algorithm

It performs $n-1$ times a relaxation of every edge in the graph.

Input: A weighted directed graph \vec{G} with n vertices and a vertex

v of \vec{G} .
Output: A label $D[u]$ for each vertex u of \vec{G} such that $D[u]$
is distance from v to u .

$D[v] \leftarrow 0$,
for each vertex $u \neq v$ of G do,
 $D[u] = +\infty$.

for $i \leftarrow 1$ to $n-1$ do,

perform the relaxation operation on (u, z) .

if $D[u] + w(u, z) < D[z]$ then

$D[z] \leftarrow D[u] + w(u, z)$;

if there are no edges left with potential relaxation operation then

return the label $D[u]$ of each vertex u .

else return " \vec{G} contains a negative-weight cycle".

Single source shortest path algo. : Synchronous Bellman-Ford.

(weighted graph ; unidirectional links).

(local variables)

int length $\leftarrow \infty$

int parent $\leftarrow -1$

set of int Neighbors \leftarrow set of neighbors

set of int { weight_{ij}, weight_{ji} | j ∈ Neighbors } \leftarrow the known values
of the incident links.

(message types)

UPDATE.

if i = i₀ then length $\leftarrow 0$,

for round = 1 to n-1 do,

send UPDATE (i, length) to all neighbors;

await UPDATE (j, length) from each j ∈ Neighbors,

for each j ∈ Neighbors do

if (length > length_j + weight_{j,i}) then

length \leftarrow length_j + weight_{j,i};

parent = j.

The source is i₀, the code is for process P_i ($1 \leq i \leq n$).

06/11/19.

Single source shortest path algorithm : Asynchronous Bellman-Ford.

(No negative weight cycles).

(local variables)

int length $\leftarrow \infty$

set of int Neighbors \leftarrow set of neighbors

set of int { weight_{ij}, weight_{ji} } \leftarrow the known values of the weights of incident lines.

(message types)

UPDATE

if $i = i_0$ then

length $\leftarrow 0$;

send UPDATE ($i_0, 0$) to all neighbors terminals.

when UPDATE UPDATE (i, length_j) arrives from j :

if ($\text{length} > (\text{length}_j + \text{weight}_{j,i})$) then

length $\leftarrow \text{length}_j + \text{weight}_{j,i}$; parent $\leftarrow j$

send UPDATE (i_0, length) to all neighbors.

Distributed file Systems .

- ✓ A distributed file system enables programs to store and access remote files exactly as they do local ones, allowing users to access files from any computer on a network.
- ✓ The performance and reliability experienced for access to files stored at a server should be comparable to that for files stored on local disks.
- ✓ Sharing of resources - key goal for distributed systems.
- ✓ Sharing of stored information - the most important aspect of distributed resource sharing.
- ✓ The design of large scale wide area read-write file storage systems poses problems of load balancing, reliability, availability and security.
- ✓ Basic of distributed file system - emulate the functionality of non-distributed file system.

* Storage file system and their properties

	Sharing	Persistence	Distributed cache/ replicas	Consistency	Example
Main Memory	x	x	x	1	RAM
File System	x	✓	x	1	Unix File System
Distributed file System	✓	✓	✓	✓	SUN NFS
Web	✓	✓	✓	x	Web Server
Distributed Shared memory	✓	x	✓	✓	
Remote objects (RMI/ORB)	✓	x	x	1	CORBA
Persistent object store	✓	✓	x	1	CORBA persistent state service
Peer-to-peer storage system	✓	✓	✓	2	

Types of consistency -

(1.: strict one copy, ✓: slightly weaker guarantee)

(2.: considerably weaker guarantee)

✓ Characteristics of the system file system.

→ responsible for the organization, storage, retrieval, naming, sharing and protection of files. They provide a programming interface.

✓ File attribute record structure

file length

creation time stamp

read timestamp

write timestamp

attribute timestamp

reference count

↓ Managed by
file system.

owner

file type

access control list ↓ user
program.

✓ File system modules (for non-distributed file system).

Directory module : relates filenames to the IDs.

File module : relates file IDs to particular files.

Access control module : checks permission for operation requested.

File access module : read or write the file data on attributes.

Block module : accesses and allocates disk blocks.

Device modules : performs disk I/O and buffering.

Device modules : performs disk I/O and buffering.

✓ Distributed file system requirements

- Initially access transparency and location transparency
- Performance, scalability, concurrency control, fault tolerance and security.
- Transparency: The design must balance the flexibility and scalability that derive from transparency against software complexity and performance.
 - Access transparency. Client programs should be unaware of the distribution of files.
 - location transparency. Client programs should see a uniform file space name space.
 - Mobility transparency. Neither client programs nor system administration tables in client nodes need not be changed when files are moved.
 - Performance transparency. Client programs should continue to perform satisfactorily by while the load on the service varies within a specified range.
 - Scaling transparency. The service can be expanded by incremental growth to deal with a wide range of loads and network sizes.