



# Synchronisation Problems

```
P():
if (value == false) {
    add myself to queue
    and block;
}
value = false;
```

(if blocks, will context switch to some other process)

```
V():
value = true;
if (queue is not empty) {
    wake up one arbitrary
    process on the queue
}
```

Executed atomically (e.g., interrupt disabled)

Does not use busy waiting. Internally has a boolean value that is set to true. And a queue of blocked processes that are initially empty (may not be FIFO). Possible to execute P() and V() out of order because they are kindof independent. Possible for a process to be stuck forever if there are other processes continually being added to the queue.

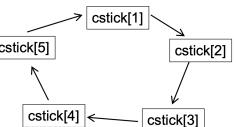
**Semantics.** Only 1 process is waken up, likely to be arbitrary but depends on implementation of queue.

**Dining Philosopher.** 5 semaphores for chopsticks and philosophers are processes. To eat, must pick up left and right. When done, put down left and right. Neighbors cannot use the chopsticks while you're using them. Ensures that no 2 process fight with each other and proper use of resource. Semaphores are chopsticks.

philosopher 1	philosopher 2	philosopher 3	philosopher 4	philosopher 5
cstick[1].P();	cstick[2].P();	cstick[3].P();	cstick[4].P();	cstick[5].P();
cstick[2].P();	cstick[3].P();	cstick[4].P();	cstick[5].P();	cstick[1].P();

- It is possible for all processes to block

Deadlock !



Avoid cycles to prevent deadlock. Or have total ordering.

## Monitor.

```
synchronized (object) {
    ...
    ...
    enterMonitor();
    ...
    exitMonitor();
    If monitor-queue is not
    empty, pick one arbitrary
    process from monitor-
    queue and unblock it
}
```

2 queues for blocked processes, monitor queue and wait queue.

**Object.wait.** Add myself to wait-queue and exit

monitor and then block myself. Done atomically. **Object.notify.** If wait queue not empty, unblock 1 process from wait-queue.

**Object.notifyAll.** If wait-queue not empty, unblock all processes in wait-queue.

**Hoare-style Monitor.** If wait-queue is non-empty, unblock all processes in wait-queue.

**Java-style Monitor.** The process that called the notify continues execution. Placing the while loop means that it will keep checking the condition after it gets unblocked.

**Conditional vs Loops.** In monitors, a while loop will be checked again when that process is allowed to continue execution. For a conditional, it will fall through. The while loop has some nice guarantees in that it could allow a particular variable to hold some value throughout the execution.

**Nested Monitors.** Calling the wait on the inner monitor releases the lock on the inner object only. This means that wait() only releases the immediate lock.

```
synchronized (ObjA) {
    synchronized (ObjB) {
        ObjB.wait();
    }
}

synchronized (ObjA) {
    synchronized (ObjB) {
        ObjB.notify();
    }
}
```

- Here Java only releases the monitor lock on ObjB and not the monitor lock on ObjA.
- Hence this piece of code will block and will not reach ObjB.notify() – deadlock!
- Different monitor implementations may differ in how nested monitors are implemented – check the spec to tell...

**Producer-Consumer.** Circular buffer with single producer and consumer. Producer put item at the end of buffer(if buffer not full) and consumer remove from the front (if buffer not empty).

object sharedBuffer;

```
void produce() {
    synchronized (sharedBuffer) {
        if (sharedBuffer is full)
            sharedBuffer.wait();
        add an item to sharedBuffer;
        if (sharedBuffer *was* empty)
            sharedBuffer.notify();
    }
}

void consume() {
    synchronized (sharedBuffer) {
        if (sharedBuffer is empty)
            sharedBuffer.wait();
        remove item from sharedBuffer;
        if (sharedBuffer *was* full)
            sharedBuffer.notify();
    }
}
```

Notification can be lost if no process is waiting.

**Reader-writer.** Multiple readers and writers accessing a file. Writer must have exclusive access but multiple readers can access the file at the same time. Flag ensures that notifications are never lost.

int numReader, numWriter; Object object;

```
void writeFile() {
    synchronized (object) {
        while (numReader > 0 || numWriter > 0)
            object.wait();
        numWriter = 1;
    }
}

void readFile() {
    synchronized (object) {
        while (numWriter > 0)
            object.wait();
        numReader++;
    }
}

// read from file;
synchronized (object) {
    numReader--;
    object.notify();
}
```

you can prove that it must be a writer who is notified

**Starvation.** To solve starvation on writers, maintain an explicit queue. Each monitor queue has at most 1 element.

```
Vector queue; // shared among all threads
int numReader, int numWriter;

// code for writer entry
Writer w = new Writer(mynname);
synchronized (queue) {
    if (numReader > 0 || numWriter > 0) {
        w.okToGo = false;
        queue.add(w);
    } else {
        w.okToGo = true;
        numWriter++;
    }
}

synchronized (w) { if (!w.okToGo) w.wait(); }
```

```
// code for writer exit
synchronized (queue) {
    numWriter--;
    if (queue is not empty) {
        remove a single writer or a batch of readers from queue
        for each request removed do {
            numberWriter++ or numberReader++;
        synchronized (request) {
            request.okToGo = true;
            request.notify();
        }
    }
}
```

```
// code for reader entry
Reader r = new Reader(mynname);
synchronized (queue) {
    if ((numWriter > 0) || !queue.isEmpty()) {
        r.okToGo = false;
        queue.add(r);
    } else {
        r.okToGo = true;
        numReader++;
    }
}

synchronized (r) { if (!r.okToGo) r.wait(); }
```

```
// code for reader exit
synchronized (queue) {
    numReader--;
    if (numReader > 0) exit(); // I am not the last reader
    if (queue is not empty) {
        remove a single writer or a batch of readers from queue
        for each request removed do {
            numWriter++ or numReader++;
        synchronized (request) {
            request.okToGo = true;
            request.notify();
        }
    }
}
```

**Barber-shop problem.** Flag ensures that notifications are never lost.

## Customer

```
Vector customQueue; // shared data among all threads

synchronized (numberChair) {
    if (numberChair > 0) numberChair--;
    else return; // leave // this releases monitor lock
}

// middle part (see next slide)

synchronized (numChair) {
    numberChair++;
}
```

## Customer – Middle part

```
// middle part
Customer myself = new Customer(mynname);
myself.done = false;
synchronized (customQueue) {
    customQueue.add(myself); // can also simulate chair here
    customQueue.notify();
}

synchronized (myself) {
    if (!myself.done) myself.wait(); // no need to use while
}

while (true) {
    Customer current;
    synchronized (customerQueue) {
        if (customerQueue.isEmpty()) customerQueue.wait();
        // no need to use "while" here because only one barber
        current = customerQueue.removeFirst();
    }

    // hair cut the current customer

    synchronized (current) {
        current.done = true;
        // this flag helps to avoid needing nested monitor
        current.notify();
    }
}
```

**PQR problem.** Notify call meant for R. This problem will be solved if we replace while loop that checks numR == numP + numQ with an if instead since only 1 item is waiting.

```

int numP, numQ, numR; Object obj;
thread1:
    while (true) {
        print P;
        synchronized (obj) { numP++; obj.notify(); }
    }
thread2: similar as above
thread3:
    while (true) {
        synchronized (obj) {
            while (numR == numberP + numberQ) obj.wait();
        }
        print R;
        synchronized (obj) { numR++; }
    }
}

```

## Consistency

**Operation:** A single invocation/response pair of a single method of a single shared object by a process

- e being an operation
  - proc(e): The invoking process
  - obj(e): The object
  - inv(e): Invocation event (start time)
  - resp(e): Reply event (finish time)
- wall clock time /  
physical time /  
real time

**Consistency.** Specifies what behavior is allowed when a shared object is accessed by multiple processes.

**History.** Sequence of invocations and responses ordered by wall clock time. For any invocation in H, corresponding response is required to be in H. Each execution of a parallel system corresponds to a history and vice versa.

**Sequential History.** Any invocation is always immediately followed by its response. If there is interleaving, then it's called concurrent.

Sequential:

inv(p, read, X) resp(p, read, X, 0) inv(q, write, X, 1) resp(q, write, X, OK)

concurrent:

inv(p, read, X) inv(q, write, X, 1) resp(p, read, X, 0) resp(q, write, X, OK)

**Legal Sequential History.** All responses satisfies the sequential semantics of the data type. Possible for sequential history not to be legal.

**Sequential Semantics.** Semantics you would get if there is only 1 process accessing that data type.

**Subhistory of a process.**  $H|p$  Subsequence of all events of the process. Thus a process subhistory is always sequential.

**Subhistory of an object.**  $H|o$  Subsequence of all events of object.

**Equivalence.** 2 histories are equivalent if they have exactly the same set of events. Same events imply **all responses are the same**. Ordering can be different as the user only cares about responses.

**Process Order.** Partial order among all events. For any 2 events of the same **process**, process order is the same as execution order. No other additional orderings.

**Sequential Consistency.** A history that is equivalent to some legal sequential history that preserves process order. Focus on results only and require that program order be preserved.

**External Order.**  $o_1 < o_2$  iff response of  $o_1$  appears

in H before invocation of  $o_2$  ( $o_1$  finishes before  $o_2$  starts). Preserve external order implies preserving program order.

**Linearisability Definition 1.** Execution equivalent to some execution st each operation happens instantaneously at some point (linearisation point) between invocation and response.

**Linearisability Definition 2.** Equivalent to some legal sequential history S and S preserves external order in H. Partial order induced by H is a subset of partial order induced by S. In other words, if a response was before an invocation in the original history, it must still precede it in S.

**Properties.** Linearisable history must be sequentially consistent. But not necessarily the other way. Both definitions are the same.

**Local Property.** H is linearisable iff for any object x,  $H|x$  is linearisable. Sequential consistency is not a local property.

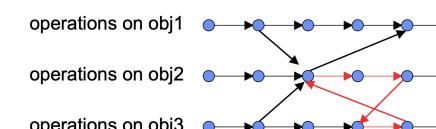
**Proof.** Using definition 2. Construct a directed graph. Directed from  $o_1$  to  $o_2$  if either is satisfied:

- $o_1$  and  $o_2$  is on the same obj x and  $o_1$  is before  $o_2$  when linearising  $H|x$  ( $o_1 \rightarrow o_2$  due to object)
- $o_1 < o_2$  in external order (response of  $o_1$  appears before invocation of  $o_2$  in H), then  $o_1 \rightarrow o_2$  due to H

Resulting graph is acyclic. So any topological sort of the graph gives a legal sequential history and the edges in the graph is a subset of the total order induced by S.

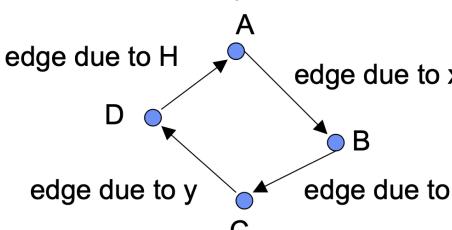
**Legal + Legal = Legal.** If  $H|x$  and  $H|y$  is linearisable and legal for  $x, y$  respectively, interleaving the operations will not affect any behavior. Since all events in S are from H (including all invocation events and parameters). So values are already captured in the params of the invocation events.

**Proof acyclic graph.** Prove by contradiction.



- To generalize, any cycle must be composed of:
  - Edges due to some object x, followed by
  - Edges due to H, followed by
  - Edges due to some object y, followed by
  - Edges due to H, followed by

Then  $H|x$  is equivalent to a serial history S where S is a total order (we can directly draw an edge from start to end), similarly for partial order induced by H, because of transitivity.



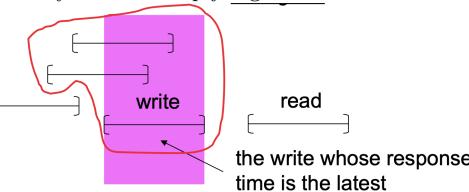
Then try to recreate the history from the cycle. From

the image, it's impossible for A to start before B finishes. The reason is that if B finishes before A, then B is before A in external order, which contradicts the edges from A to B. Contradiction here arises because we have to put the start of C to be before the end of D which is not possible.

**Register.** A kind of abstract data type - a single value that can be read and written.

**Atomic register.** Implementation ensures linearisability of history.

**Regular register.** When a read does not overlap with any write, the read returns the value written by 1 of the most recent writes. When a read overlaps, with  $\geq 1$ , read returns value written by 1 of the most recent writes/1 of the overlapping writes. Atomic registers must be regular. **Does not imply sequential consistency**, similarly sequential consistency does not imply regular.



**Safe registers.** When a read does not overlap with any write, then it returns the value written by 1 of the most recent writes. When a read overlaps with  $\geq 1$  writes, return anything. Regular registers are safe. Does not imply sequential consistency and vice versa.

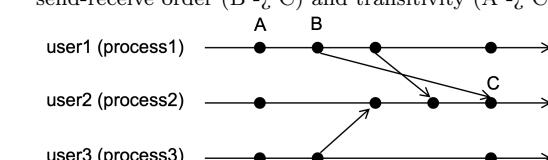
## Clocks

Process can receive 3 kinds of atomic actions/events. Local computation, send/receive single message to/from process.

**Communication Model.** Point-to-point. Error free, infinite buffer and potentially out of order.

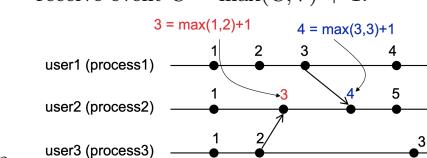
**Software Clocks.** Capture event ordering that are visible to users who don't have physical clocks.

**Visible ordering.** Process order ( $A \prec B$ ), send-receive order ( $B \prec C$ ) and transitivity ( $A \prec C$ )

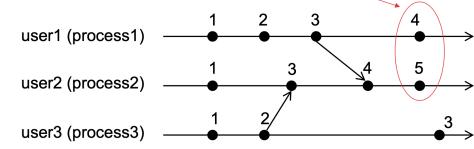


**Happened-before relation.** Captures ordering that is visible to users when there is no physical clock. Creates a partial order among events and preserves process order, send-receive order and transitivity.

**Logical Clock.** Each event has a single integer as its logical clock value. Each process has a lock counter C. Increment C at each local computation and send event. When sending a message, logical clock value V is attached to the message. At each receive event  $C = \max(C, V) + 1$ .



**Properties.** If event S happens before T, then the logical clock value of S is smaller than the logical clock value of T. The reverse may not be true.



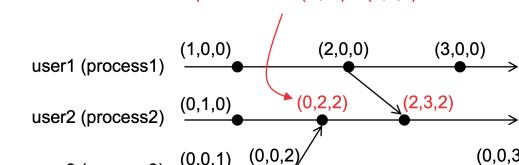
For total order, extend with process number.  $s.c$  denotes the value of c in state s,  $s.p$  indicates the process it belongs to. Timestamp of any event is a tuple  $(s.c, s.p)$ . Then  $(s.c, s.p) < (t.c, t.p) \Leftrightarrow (s.c < t.c) \vee ((s.c = t.c) \wedge (s.p < t.p))$

**Vector clock.** Event s happens before event t  $\Leftrightarrow$  the vector clock value of s is smaller than that of t.

**Partial order on clock value.** Each event has a vector of n integers as its vector clock value.  $v1 = v2$  if all the fields are the same.  $v1 \leq v2$  if every fields in  $v1 \leq$  to corresponding field in  $v2$ .  $v1 < v2$  if  $v1 \leq v2$  and  $v1 \neq v2$ . Partial order counter example:  $[2, 0, 1], [1, 2, 0]$ .

**Protocol.** Each process i has a local vector C. Increment  $C[i]$  at each local computation and send event. When sending a message, vector clock V is attached to the message. At each receive event,  $C = \text{pairwise-max}(C, V)$ . Then increase  $C[i]$  by 1.

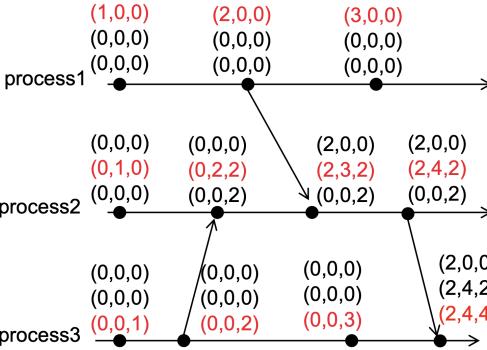
$$C = (0, 1, 0), V = (0, 0, 2) \\ \text{pairwise-max}(C, V) = (0, 1, 2)$$



**Properties Proof.** Event s happened before t then vector clock value of s  $\leq$  vector clock value of t. Enumerate all possible cases. Smaller than relation among vector clock values is transitive. Vector clock value of s  $\leq$  vector clock value of t implies s happened before t. If s and t on the same process, done. If s is on p and t is on q, let VS be s vector clock and VT be t.  $VS < VT \Rightarrow VS[p] \leq VT[p]$  Must be a sequence of message from p to q after s and before t.

**Matrix Clock Protocol.** Each event has n vector clocks, 1 for each process.  $C[i]$  is the process  $i$  principle vector, where  $C$  is a matrix. Principle vector is the same as vector clock before. Non-principle vectors are just piggybacked on messages to update knowledge.

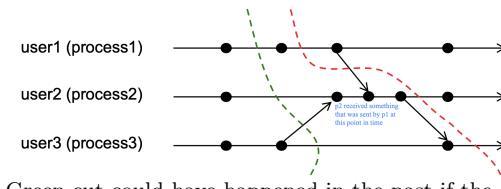
**Protocol.** Increment  $C[i]$  at each local computation and send event. When sending a message, all n vectors are attached to the message. At each receive event, let V be the principle vector of the sender.  $C = \text{pairwise-max}(C, V)$ . Then increase  $C[i]$  by 1. For non-principle vector C on process i, suppose it corresponds to process j. At each receive event, let V be the vector corresponding to process j as in the received message.  $C = \text{pairwise-max}(C, V)$ .



## Global Snapshot

A snapshot of local states of  $n$  processes such that the global snapshot could have happened sometime in the past (users cannot tell the difference).

**Consistent Snapshot.** Snapshot of local states on n processes st global snapshot could have happened sometime in the past.

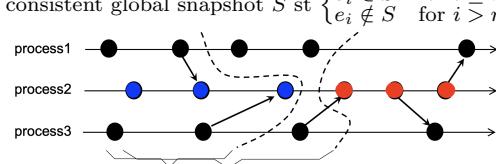


Green cut could have happened in the past if the processes had made progress faster/slower but not the red cut.

**Global Snapshot.** A set of events st if  $e_2$  is in the set and  $e_1$  is before  $e_2$  in **process order**, then  $e_1$  must be in the set.

**Consistent Global Snapshot.** A global snapshot st if  $e_2$  is in the set and  $e_1$  is before  $e_2$  in send-receive order, then  $e_1$  must be in the set. Transitive relation implied in global snapshot definition. In other words, if a receive event is in the set, corresponding send event must be in the set.

**Existence.** For any positive integer  $m$ , there is a consistent global snapshot  $S$  st  $\{e_i \in S \text{ for } i \leq m\}$   $\{e_i \notin S \text{ for } i > m\}$

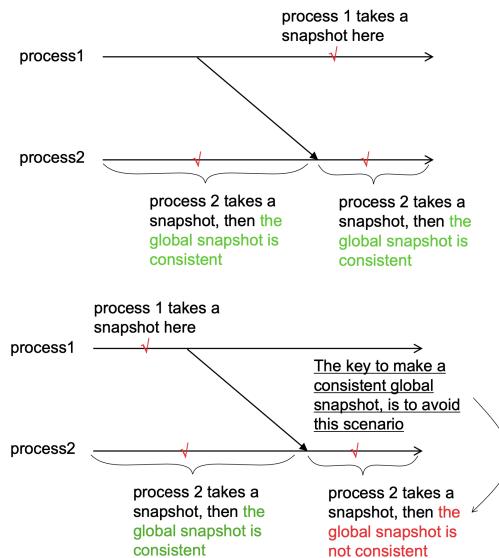


3 kinds of events:

- Blue events must be in  $S$
- Red events cannot be in  $S$
- Other events that may/may not be in  $S$

**Capturing CGS.** Look at communication model. **Communication Model.** No message loss, unidirectional and FIFO.

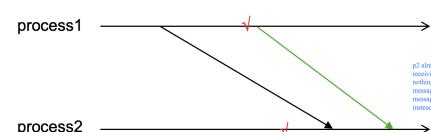
**Ensuring FIFO.** Each process maintains a message number counter for each channel and stamps each message sent. Receiver only deliver messages in order.



**Chandy Lamport.** Each process either red (has taken local SS) or white (hasn't taken). Initiated by single process by turning itself from white to red. Once a process turns red, send out market messages to all other processes. Upon receiving marker, process turns Red. There will be a total of  $n(n + 1)$  market messages. **On-the-fly.**

- M is sent **before** the local snapshot (on the sender) and received **before** the local snapshot (on the receiver): Not "on-the-fly" (why?)
- M is sent **after** the local snapshot and received **after** the local snapshot: Not "on-the-fly" (This is not CGS)
- M is sent **after** the local snapshot and received **before** the local snapshot: "on-the-fly"
- M is sent **before** the local snapshot and received **after** the local snapshot: "on-the-fly"

On the fly messages are messages that are sent before sender's local snapshot and received after receiver's local snapshot.



Process2 records all the messages received in the window from its local snapshot until it received the marker message. **Property.** If G and H are both CGS, then  $G \cap H$  and  $G \cup H$  are also CGS.

**Proving tip.** To show CGS, need to show that its a global snapshot (process order definition) then show that it is in CGS (send-receive order).

## Message Ordering

**Causal Order.** If  $s_1$  happened before  $s_2$ , and  $r_1$  and  $r_2$  are on the same process, then  $r_1$  must be before  $r_2$ . **Protocol.** Each process maintains a  $n \times n$  matrix

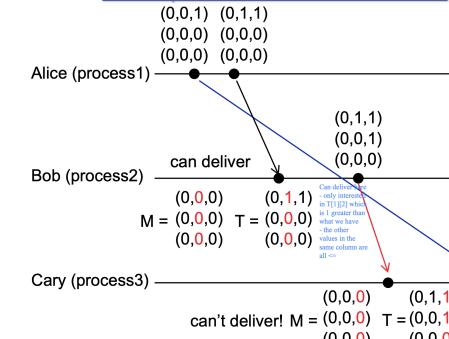
$M$ .  $M[i][j]$  is the number of messages sent from  $i$  to  $j$  as known by process  $i$ . If process  $i$  sends a message to process  $j$ , then on process  $i$ ,  $M[i][j] += 1$ . Piggyback  $M$  on the message.

**Receiver.** Local matrix on  $j$  is  $M$ , piggybacked matrix from  $i$  is  $T$ .

**Delivery conditions.** Deliver the message and set  $M$  to be pairwise max of  $M, T$  if

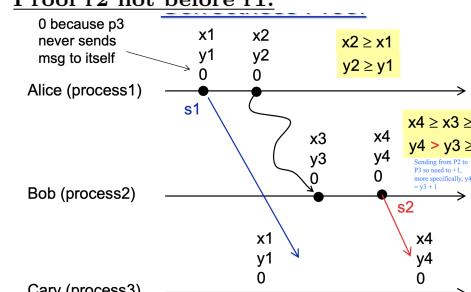
$$\begin{cases} T[k][j] \leq M[k][j] & \text{for all } k \neq i \\ T[k][j] = M[k][j] + 1 & \text{otherwise delay} \end{cases}$$

message.  $M[i][j]$  on process  $j$  takes consecutive values. If the entry value  $> 1$  larger than local, then there is another message in propagation. Column  $j$  should only have a difference of 1.



**Proof of correctness.**  $s_1, s_2, r_1, r_2$  where  $s_1$  happen before  $s_2$ . Need to prove both  $r_1$  and  $r_2$  eventually delivered and  $r_2$  not before  $r_1$ .

**Proof  $r_2$  not before  $r_1$ .**



conclusion so far:  $x_4 \geq x_1$  and  $y_4 > y_1$

**Proof by contradiction.** Assume red delivered before blue. Then after delivering red, 3rd column of  $M$  on process3 will contain  $x_4, y_4, 0$ . Since  $M$  **never decreases**, blue cannot be delivered anymore.

**Proof  $r_1$   $r_2$  eventually delivered.** At any given time, there must be 1 message that can be delivered to the process. Consider a receiver process  $j$  and its corresponding column in  $M$ ,  $x_1, \dots, x_i, \dots, x_n$

**Successor message.** Non-empty set of all undelivered messages. For each sender  $i$  in the set, there is an undelivered message whose matrix  $T$  has  $x_i + 1$  as a value in its column. Set of successor messages is non-empty, at most 1 message from any given sender and at most  $n - 1$  messages in total.

**Top successor message.** There must be a send event  $s$  that does not have any other send events that happened before  $s$ . Possible to have multiple such  $s$  events. Any top successor message can be delivered.

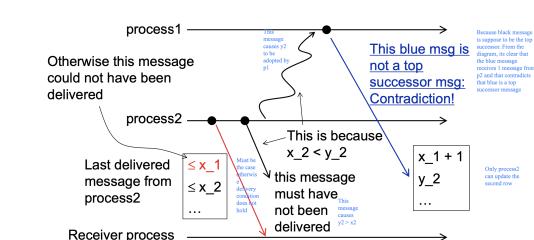
$x_1$	$x_{i+1}$
$x_2$	$y_2$
$\dots$	$\dots$
$x_n$	$y_n$

The matrix column on the receiver (process  $j$ )

The matrix column in the message

We prove the case for  $i = 2$ . Other cases are the same.

Proof by contradiction. Suppose  $x_2 < y_2$ .

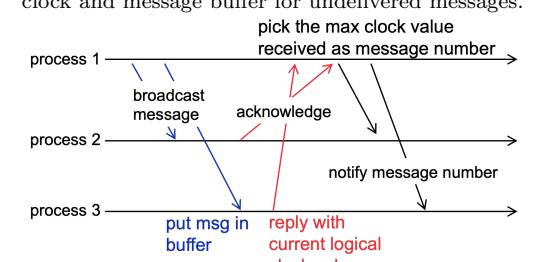


**Broadcast.** Every message sent to everyone including itself. Modeled as  $n$  point-to-point messages.

**Total Ordering.** All messages delivered to all process in exactly the same order (atomic broadcast). Only applies to broadcast messages. Total order does not imply causal order and vice versa.

**Coordinator.** Special process. To broadcast message, send message to coordinator. Coordinator assigns sequence number to message and forwards the message to all processes with the sequence number. Messages delivered according to sequence number order. Problem is that coordinator has too much control.

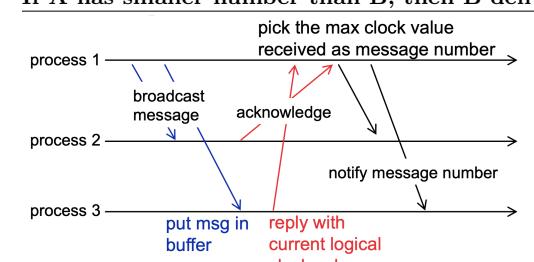
**Skeens Algorithm.** Each process maintain logical clock and message buffer for undelivered messages.



**Buffer messages.** Delivered/removed if all messages in the buffer have been assigned numbers or this message has the smallest number.

**Correctness proof.** All messages will be assigned message numbers and delivered.

**If A has smaller number than B, then B delivered**



Suppose A delivered on process3 after B. Then A must have been placed in buffer after B was delivered. This means that A must have a larger number than B - contradiction.

# Leader Election

Trivially solves mutual exclusion and total order broadcast. **Anonymous ring.** No unique identifiers, node can only send messages to its neighbors. Impossible using deterministic algorithms. Same initial state, algo, step and final state.

**Chang-Roberts Algorithm.** Each node have unique id. Nodes only send messages clockwise. Each node acts on its own.

**Protocol.** Node send election message with its own id clockwise. Election message is forwarded if id in messages larger than own message. Otherwise message is discarded. Node becomes leader if it sees its own election message.

**Performance.** For distributed systems, communication is bottleneck. Performance described as message complexity. With  $n$  nodes on the ring, whats the total number of messages incurred by protocol.

**Best/Worse case.** Best case:  $2n - 1$ . Worse cast:  $\frac{n(n+1)}{2}$

**Average performance.**  $O(n \log n)$  Average is taken over all possible ordering of the nodes on the ring and each ordering has the same probability. Total ordering excluding isomorphic ones  $(n-1)!$ .  $x_k$  is the number of messages caused by election message from node  $k$ . Want to find  $E[\sum x_k]$  for  $k$  summed from 1 to  $n$ . Then by linearity of expectation,

$$E[\sum x_k] = \sum E[x_k]$$

**Probability trick.**  $Pr[x=1] = \frac{n-k}{n-1}$ . Since  $x_k$  is the number of messages caused by node  $k$ , total  $n-1$  nodes other than  $k$ , of which  $n-k$  nodes have id larger than  $k$ . Then  $Pr[x_k = 2|x_k > 1] = \frac{n-k}{n-2}$ .

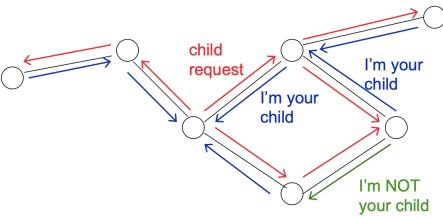
Simplify  $p = \frac{n-k}{n-1}$  and upper bound everything by  $p$ . Let  $y$  be rv denoting number of tickets we buy until we win for the first time, assuming probability of winning is  $q$ , then  $Pr[y = (i+1)|y > i] = q$ . Then  $E[y] = \frac{1}{q}$ . So  $E[x_k] \leq E[y] = \frac{1}{q} = \frac{n-1}{n-k}$

$$\begin{aligned} \sum_{k=1}^n E[x_k] &= n + \sum_{k=1}^{n-1} E[x_k] \\ &< n + \sum_{k=1}^{n-1} \frac{n-1}{n-k} \\ &= n + (n-1) \sum_{k=1}^{n-1} \frac{1}{k} \\ &= n + (n-1)O(\log n) = O(n \log n) \end{aligned}$$

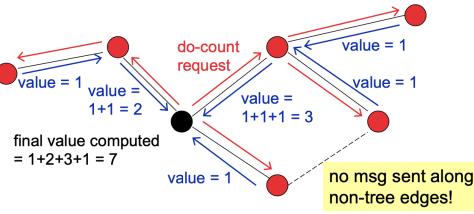
**General graph (known  $n$ ).** Complete graph or any connected graph.

**Complete graph.** Each node sends its id to all other nodes. Wait until you receive  $n$  ids. Biggest id wins.

**Any connected graph.** Flood your id to all other nodes. Wait until you receive  $n$  ids. Biggest id wins.  
**Spanning Tree Construction to count  $n$ .** No centralised coordinator, possible for multiple nodes to initiate this. Goal of the protocol is that each node knows its parent and children. Each node sends child request. If already a child of another node, reject request.



**Counting nodes.** Initiator sends a do-count request. Recursively, children will respond with  $1 +$  number of children they have.



**Spanning tree purpose.** Broadcast (root doesn't need to send out  $n$  messages), aggregations, calculating maxmin.

**Impossibility for unknown ring size.** Modelling a randomised algorithm (a deterministic algorithm with 1 random input).

**Proof.** Assume there is an algorithm to solve the problem. Terminate with probability 1 implies that there at least exists some bit string input that terminates after  $r$  (constant) phases in ring 1. In ring 2, same bit string may be input for P1 and P2. Ring 1 and 2 indistinguishable for P1, so P1 declares leader after  $r$  phases in ring 2 as well. Similarly for P2. In ring 2, there will be 2 leaders, contradiction.



**For known ring size.** Each node has phase number initialised to 1. Each message sent has phase number attached.

**Phase 1.** Each node pick random id and run Chang-roberts. Winner proceed to next phase, losers only forward messages in future phases. A node know its the winner if it sees its own id after exactly  $n$  hops (ring size). Stop if there is a single winner. By definition there will be a single winner if protocol stops.

**Proof of termination.** Define the  $i$ th good phase to be the phase that results in  $i$ th loser. Need exactly  $n-1$  good phase. Let  $r$  be the phase number.  $x_i$  be the number of phases between the  $i$ th and  $i+1$ th good phase. As  $r \rightarrow \infty$ ,  $Pr[x_i > \frac{r}{n-1}] \rightarrow 0$

## Distributed Consensus

Set of nodes want to agree on something. Each node has an input.

**Version 0.** No failures. Everyone send its input to everyone else. Take the majority of the  $n$  values (with tie-breaking favoring 0). Decide and done.

**Failure Model.** Nodes can experience crash failure. Communication channels are reliable.

**Synchronous timing model.** Message delay has known upper bound  $x$  and node processing delay has a known upper bound  $y$ . Implies possibility of accurate failure detection.

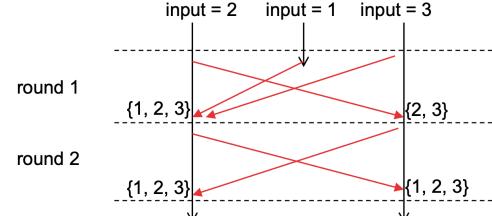
**Goals.** **Termination** - all nodes that have not failed eventually decides. **Agreement** - all nodes that decide should decide on the same value. **Validity** If all nodes have the same initial input, that value should be the only possible decision value. Otherwise nodes can decide on anything (except they still need to satisfy Agreement).

**Version 1.** Node crash failures and synchronous timing model.

**Rounds.** Under synchronous timing model, all processes proceed in inter-locked rounds. In each round, each process does some local computation, sends 1 message to every other process and receives 1 message from every other process.

**Round implementation.** Message propagation delay + local processing delay + clock error. Assume each process has a physical clock with some bounded clock error. Message sent in a round must be received by the end of that round on the receiver (may be earlier than that). Each message has a round number attached on it.

**Protocol.** Each process sends its input to all the others and pick min/max. If there are failures, process may send input to only a subset of the processes. Keep forwarding the values until all processes have received.



Need 1 extra round for each failure ( $f+1$  rounds to tolerate  $f$  failures).  $f$  is input to protocol. With  $f$  crash failures, any consensus protocol will take  $\Omega(f)$  rounds (lower bound).

**Proof.** Termination is obvious since there is no wait. Validity obvious since S will only have a single element.

**Agreement proof.** Prove that S is the same on all nodes when protocol terminates.

**Nonfaulty.** Nodes in round  $r$  that did not crash by the end of round  $r$ .

**Good round.** No node failure during that round. With  $f+1$  rounds and  $f$  failures, at least 1 good round.

**Claims.**

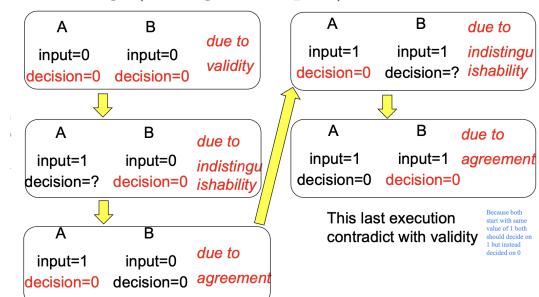
- At end of any good round  $r$ , all non-faulty nodes during round  $r$  have the same S.
- Suppose  $r$  is a good round, and consider any round  $t$  after round  $r$ . During round  $t$ , value of S on any non-faulty nodes doesn't change.
- All non faulty processes at round  $f+1$  have the same S.

**Coordinated attack problem.** Failure model nodes dont fail but communication channels might (may drop unbounded number of messages). Timing model is synchronous.

**Achievability.** Impossible to reach goal using

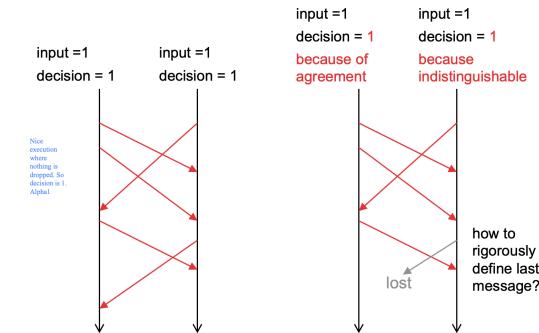
deterministic algorithm since communication channel can drop all messages.

**Indistinguishable.** 2 executions are indistinguishable for some node if the node sees the same things (message and inputs) in both executions.



Contradiction arises because with same inputs, should decide on 1 but decided on 0 instead. From the first execution to the next, channel drops all messages from B to A. So A cannot make a decision. Since B has the same things for the 2 executions, become indistinguishable.

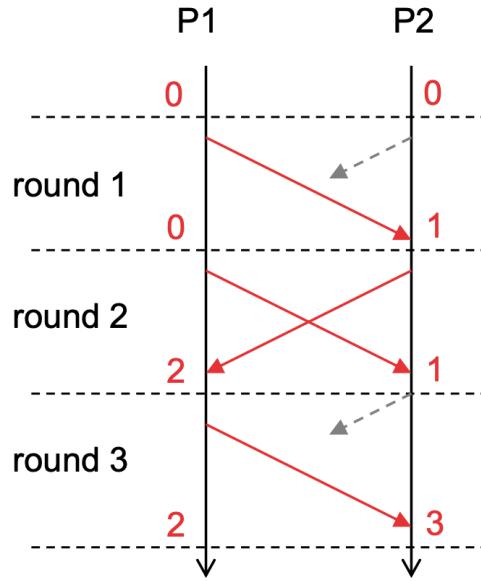
**Weakened validity.** Assuming input either 0/1. If all nodes start with 0, decision should be 0. If all nodes start with 1 and no message lost throughout execution decision should be 1. Otherwise allowed to decide on anything. Channel does not drop messages when all nodes start with 1. Still impossible using deterministic algorithm.



On the left is the nice execution where nothing is dropped, so decision is 1. On the right, if we drop the last message(sent time), there is no way of either process knowing. Indistinguishable for the right process since everything is the same. Keep repeating by dropping the last message for each execution until there are no more messages.

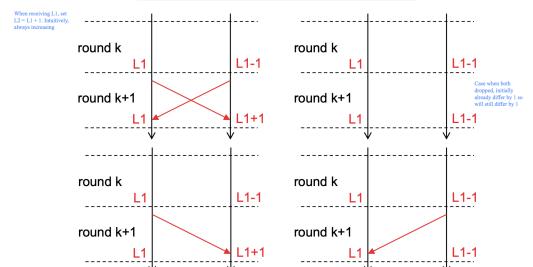
**Adversary.** Process can flip coins and adversary try to maximise  $P_e$  by setting inputs of process, causing message loss. Outcome of random flips unknown to adversary.

**Randomised Algorithm.** Consider 2 processes. Algorithm has predetermined number of rounds,  $r$ . Adversary determines which messages get lost before seeing random choice. P1 picks a random integer  $bar \in [1, \dots, r]$  at the start. Protocol allows P1 and P2 to each maintain a level variable (L1, L2) st level is influenced by adversary but L1, L2 differ by at most 1.



**bar**, input and current level attached on all messages. When receiving, update own level variable to be  $level_{self} = level_{received} + 1$ . **L1/L2** never decreases and differ by at most 1.

#### Inductive Proof for the Lemma



**Decision rule.** At the end of  $r$  rounds, P1 choose 1 iff P1 knows P1 and P2 inputs are 1 and  $L1 \geq bar$ . Implies that P1 decide on 0 if doesn't see P2 input. P2 decides on 1 iff P2 knows P1 P2 inputs are both 1 and P2 knows bar and  $L2 \geq bar$ . Implies that P2 decide on 0 if it doesn't see P1 input or doesn't see bar.

**Error Occurrence.** For P1 and P2 to decide on different values, one must decide on 1 and the other decides on 0. Both inputs must be 1. All cases mutually exclusive and happens with probability  $\frac{1}{r}$  3 cases:

- P1 see P2 input but P2 dont see P1 input/bar.  $L1 = 1$  and  $L2 = 0$ . Error occurs when  $bar = 1$ .
- P2 see P1 input and bar but P1 dont see P2 input.  $L1 = 0$ ,  $L2 = 1$ . Error occurs only when  $bar = 1$ .
- P1 sees P2 input, P2 sees P1 input and bar. Then error occurs when  $bar = \max(L1, L2)$ .

**Correctness Proof.** Termination is obvious since  $r$  rounds. Validity - if all rounds start with 0, decision always 0. If all nodes start with 1 and no messages lost throughout execution, decision should be 1. Since no message lost then  $L1 = L2 = r$ . Agreement happens with probability  $1 - \frac{1}{r}$

**Consensus with Node Crash - Async.** Failure model same as v1. Asynchronous - process and message delay are finite but unbounded. Cannot define a round since if there is no message for a long time, unsure if sender failed or message delayed. Goals are still termination, agreement and validity.

**FLP Theorem.** Distributed consensus problem under async communication model is impossible to solve even with a single node crash failure. Fundamentally protocol is unable to accurately detect node failure.

**Formalism.** Goal is to abstract execution of any possible deterministic protocol. Each process has 2 variables, input  $\in \{0, 1\}$  and decision  $\in \{null, 0, 1\}$  decision is initially null and can be written exactly once.

**Communication channel.** Messages are on the fly. **Message system.** Captures state of all communication channels. Is the set containing pairs of  $(p, m)$  where  $m$  is on the fly to process  $p$ . All messages are distinct.

**Send.** add a (dest, content) to the message system.

**Receive.** When invoked by process  $p$ . Remove some  $(p, content)$  from message system and return content **OR** leave message system unchanged and return null. No fifo ordering means no guarantee on message order. Not a special case of out-of-order. FLP holds when receive is blocking. Can use blocking to implement non-blocking.

**Step.** Takes system from 1 global state to another. By executing the following on process  $p$ . Receive a message  $m$  (can be null). Based on  $p$  local state and  $m$ , send an arbitrary but finite number of messages and change  $p$  local state to a new one.

**Event.**  $(p, m)$  describes  $p$ 's receiving  $m$ . Events are input to state machine that cause state transitions. An event can be applied to global state  $G$  if either  $m$  is null or  $(p, m)$  in message system.

**Execution.** For any protocol, can be abstracted to be an **infinite** sequence of events. Each execution may be different. Can always make a protocol not to terminate. Each process must be able to handle null messages. Decisions are made when decision variable is set. Necessary to properly defined failed processes.

**Schedule.**  $\sigma$  is a sequence of events that captures execution of some protocol.  $\sigma$  can be applied to  $G$  if events can be applied to  $G$  in the order in  $\sigma$ .

$G' = \sigma(G)$  means if we apply  $\sigma$  to  $G$  we will end up with  $G'$ .  $\sigma(G)$  may or may not be applied to  $G$ .

Schedule creates a total order on events.

**Reachable.** Given a consensus protocol  $A$ , a global state  $G_2$  is reachable from  $G_1$  if there is a schedule  $\sigma$  of  $A$  such that  $G_2 = \sigma(G_1)$ .

**Async Systems.** **Processes have unbounded but finite delay.** Nonfaulty process takes infinite number of steps. A faulty one takes a finite number of steps. If we consider only finite sequences, cannot distinguish faulty and nonfaulty processes. **Messages have unbounded but finite delay.** Every message is eventually delivered. If there is a message  $(p, m)$  in the message system and  $p$  invokes receive multiple times, then the message system can only return null finite number of times. At most 1 faulty process.

**Proof technique for FLP.** Act as adversary to defeat consensus protocol. Scheduler can pick which messages to deliver and which process will take the next step (under constraints of async system). Goal is to prevent protocol from deciding.

**Global states.**  $G$  is **0-valent** if 0 is the only possible decision reachable from  $G$ . Processes in  $G$

may or may not yet decided on 0, but otherwise, they will eventually decide on 0.  $G$  is **1-valent** if 1 is the only possible decision reachable from  $G$ .  $G$  is **univalent** if  $G$  is either 0-valent or 1-valent.  $G$  is **bivalent** if its 0-valent and 1-valent.

**Lemmas.**

• For any protocol  $A$ , there exists a bivalent initial state. Prove by contradiction.

• Let  $\sigma_1, \sigma_2$  be 2 schedules st the set of processes executing steps in  $\sigma_1$  are disjoint from the set that executes  $\sigma_2$ . Then for any  $G$  that  $\sigma_1, \sigma_2$  can both be applied, we have  $\sigma_1(\sigma_2(G)) = \sigma_2(\sigma_1(G))$ . Prove by induction.

• Let  $G$  be a global state, and  $e = (p, m)$  is an event that can be applied to  $G$ . Let  $W$  be the set of global states that is reachable from  $G$  without applying  $e$ , then  $e$  can be applied to any state in  $W$ .

• Let  $G$  be a bivalent state, and  $e = (p, m)$  is any event that can be applied to  $G$ . Let  $W$  be the set of global states that is reachable from  $G$  ( $G$  is in  $W$ ) without applying  $e$ , and  $V = e(W)$  to be the set of global states by applied  $e$  to the states in  $W$ . Then  $V$  contains a bivalent state.

**n processes; at most f failures; f+1 phases; each phase has two rounds**

**Code for Process i:**

```

V[1..n] = 0; V[i] = my input;
for (k = 1; k <= f+1; k++) { // (f+1) phases
    send V[i] to all processes;
    set V[1..n] to be the n values received;
    if (value x occurs (> n/2) times in V) proposal = x;
    else proposal = 0;
}

round for all-to-all broadcast { if (k==i) send proposal to all; // I am coordinator
    receive coordinatorProposal from the coordinator
    decide whether to listen to coordinator { if (value y occurs (> n/2 + f) times in V) V[i] = y;
        else V[i] = coordinatorProposal;
    }
    decide on V[i];
}

```

if the coordinator, and my proposal to everyone is the same as the previous step. So everyone receives the proposal from the previous step.

**Proof.** **Lemma 1:** If all good processes  $P_i$  have  $V[i] = y$  at the start of phase  $k$ , then this remains true at the end of phase  $k$ . There are  $n - f$  good processes. This means that  $n - f > \frac{n}{2} + f$ , so they form majority. **Lemma 2:** If coordinator in phase  $k$  is good, then all nonfaulty processes  $P_i$  have the same  $V[i]$  at the end of phase  $k$ . **Case 1** Coordinator has proposal value as  $x$  (must be unique to coordinator). On coordinator,  $x$  appears  $> \frac{n}{2}$  times in  $V$  means that  $> \frac{n}{2} - f$  must be from good processes. On other processes,  $x$  appears  $> \frac{n}{2} - f$  times means its impossible for  $y$  to appear  $> \frac{n}{2} + f$  times in  $V$ . **Case 2** Coordinator has proposal value as 0. On coordinator, no value  $x$  appears  $> \frac{n}{2}$  times in  $V$ . So on other processes its impossible for  $y$  to appear  $> \frac{n}{2} + f$  times in  $V$ .

**Correctness.** **Termination** obvious since  $f + 1$  phases. **Validity** follows lemma 1. **Agreement** - with  $f + 1$  phases, at least 1 is a deciding phase. Then by lemma 2, after deciding phase, all good process have same  $V[i]$ . From lemma 1, in following phases,  $V[i]$  on good processes dont change.

Failure Model and Timing Model	Consensus Protocol
Ver 0: No node or link failures	Trivial – all-to-all broadcast
Ver 1: Node crash failures; Channels are reliable; Synchronous;	(f+1)-round protocol can tolerate f crash failures
Ver 2: No node failures; Channels may drop messages (the coordinated attack problem)	Impossible without error Randomized algorithm with 1/r error prob
Ver 3: Node crash failures; Channels are reliable; Asynchronous;	Impossible (the FLP theorem)
Ver 4: Node Byzantine failures; Channels are reliable; Synchronous; (the Byzantine Generals problem)	If n <= 3f, impossible. If n >= 4f + 1, we have a (2f+2)-round protocol.
	How about 3f+1 < n < 4f ? Solvable

## Self-stabilisation

State of distributed system is either legal/illegal. Definition based on application semantics. Code on each process assumed to be correct all the time.

**Self-stabilising.** Starting from any (legal/illegal) state, protocol eventually reach legal state if there are no faults. Once the system is in a legal state, it will only transit to other legal states unless there are faults. Intuitively will always recover from faults and once recovered, will remain recovered forever.

**Self-stabilising Spanning Tree.** Given  $n$  processes connected by undirected graph and 1 special process  $p_1$ , construct a spanning tree rooted at  $p_1$ . Not all

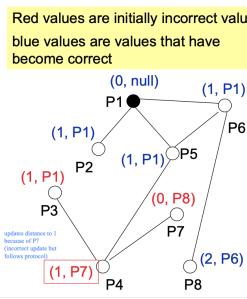
processes can communicate with all processes directly. Can be used to compute shortest path.

**Algorithm.** Each process maintains 2 variables, parent (whose my parent) and dist (distance to root, never negative). Variables are continuously updated. At any given point of time, values of 2 variables can be wrong due to faults.

On P1 (executed periodically):  
 - dist = 0; parent = null;

On all other processes (executed periodically):

- Retrieve dist from all neighbors
- Set my own dist = 1 + (the smallest dist received)
- Set my own parent = my neighbor with the smallest dist (with tie breaking if needed)



updates distance to 1 because of P1's update but ignores P2's update

**Comparison with leader election.** Leader election builds spanning tree from scratch. Self-stabilising either builds from scratch or incrementally repairs.

**Formality.** A phase to be the min time period where each process has executed its code at least once. A level is  $A_i$  be the length of shortest path from process i to root.  $dist_i$  be value of dist on process i. When topology does not change,  $A_i$  fixed but  $dist_i$  may change.

**Level properties.** A node at level X has at least 1 neighbor in level X-1. A node at level X can only have neighbors in level X-1, X and X+1.

**Correctness proof.** **Lemma 1:** at end of phase 1,  $dist_1 = 0$  and  $dist_i \geq 1, \forall i \geq 2$ . **Lemma 2:** At end of phase r, any process i whose  $A_i \leq r - 1$  will have  $dist_i = A_i$  and any process i whose  $A_i \geq r$  will have  $dist_i \geq r$

**Technique for self-stabilisation proof.** Process may take multiple actions in a phase and actions may be done in parallel - cannot assume a serialisation of all actions. Steps to proof:

- Prove that the t actions will not roll back what is already achieved so far by phase r (no backward move)
- Prove that at some point, each node will achieve more (forward move)
- Prove that the t actions will not roll back the effects of the forward move after the forward move happens (no backward move after forward move)

**Claim 1:** The blue conditions hold *throughout* phase r+1.

Proof: Use an induction on t. If t = 0, trivial. Next assume the statement holds for t and consider action (t+1) by some node i. (Cannot assume that action (t+1) happens after the t actions.)

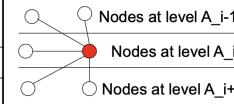
for nodes with	already know: phase r	want to show: phase r+1
$A_i \leq r - 1$	$dist_i = A_i$	$dist_i = A_i$
$A_i = r$	$dist_i \geq r$	$dist_i = r$
$A_i \geq r + 1$	$dist_i \geq r$	$dist_i \geq r + 1$

Proof (continued):

Case I:  $A_i \leq r - 1$ . Node i has at least one neighbor Y at level  $A_i - 1$ . By the blue condition and by our inductive hypothesis, we know that Y's dist value is exactly  $A_i - 1$ . Furthermore, node i only has neighbors at level  $A_i - 1$ ,  $A_i$ , and  $A_i + 1$ . By the blue condition and by our inductive hypothesis, the dist value on any of these neighbors will not be smaller than  $A_i - 1$ . Hence the min dist seen by node i is  $A_i - 1$ . Hence node i will set its dist value to be  $A_i$ .

Case II:  $A_i \geq r$ . Similar...

for nodes with	already know: phase r	want to show: phase r+1
$A_i \leq r - 1$	$dist_i = A_i$	$dist_i = A_i$
$A_i = r$	$dist_i \geq r$	$dist_i = A_i$
$A_i \geq r + 1$	$dist_i \geq r$	$dist_i \geq r + 1$



**Claim 2:** Each node i will satisfy the red conditions at some point of time during phase r+1.

Proof: By definition of a phase, node i takes at least one action in phase r+1. Consider three possible cases for node i:  $A_i \leq r - 1$ ,  $A_i = r$  and  $A_i \geq r$ . Also, leverage the fact the blue conditions always hold throughout phase r+1.

for nodes with	already know: phase r	want to show: phase r+1
$A_i \leq r - 1$	$dist_i = A_i$	$dist_i = A_i$
$A_i = r$	$dist_i \geq r$	$dist_i = A_i$
$A_i \geq r + 1$	$dist_i \geq r$	$dist_i \geq r + 1$



**Claim 3:** For each node after it first satisfies the red condition, it will continue to satisfy the red condition for the remainder of phase r+1.

Proof: Enumerate three cases, and also to leverage the fact that the blue conditions always hold throughout phase r+1.

for nodes with	already know: phase r	want to show: phase r+1
$A_i \leq r - 1$	$dist_i = A_i$	$dist_i = A_i$
$A_i = r$	$dist_i \geq r$	$dist_i = A_i$
$A_i \geq r + 1$	$dist_i \geq r$	$dist_i \geq r + 1$



**Theorems.** After H+1 phases,  $dist_i = A_i$  on all processes. H is the length of the shortest path from the most far away process to the root. Directly follows from earlier lemma. After H+1 phases, dist and parent values on all processes are correct. Each process has a single parent pointer except the root. So the graph has n nodes and n-1 edges. Each process has a path to the root, thus the graph is connected. Hence its a spanning tree.