

ALGORITHMIC APPROACH TO D.S.

2020/2021



20/09/21

LEADER ELECTION

COMPLETE
GRAPH LB

WORST CASE

- Process i receive IDs in increasing order
 - it receive $M-1$ msgs
 - it broadcast ($M-1$ msgs) for each one he receive
 - there are M processors
- $$\Rightarrow \# \text{MSGS} = (M-1)(M-1)M = O(M^3)$$

BEST CASE:

- Process i receives IDs in mon-increasing order
 - it receive $M-1$ msgs
 - there are M processors
- $$\Rightarrow \# \text{MSGS} = (M-1)M = O(M^2)$$

CHANG / ROBERTS

- ASSUME that exists a UNIDIRECTIONAL RING of non faulty processors
- PROCLAMATION is needed ($<\text{End}$) msg

ALGO

I_p

$\{ M = \emptyset \}$

$M = p$

SEND M

$\#P$ is Processor 10

M_p

$\{ \text{RECEIVING } J \}$

IF ($J > M$)

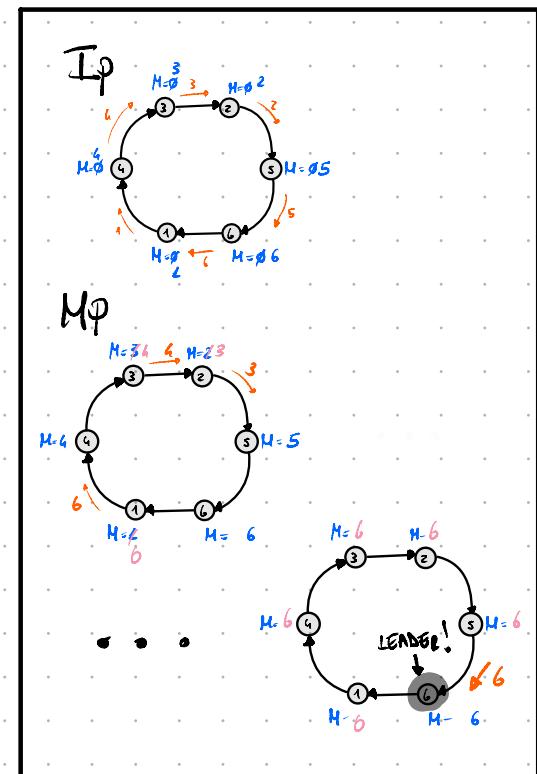
$M = J$

SEND M

IF ($M = p$)

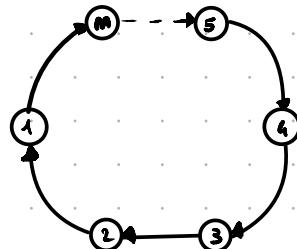
PROCLAMATION

STOP



NOTE IGNORE PROCLAMATION

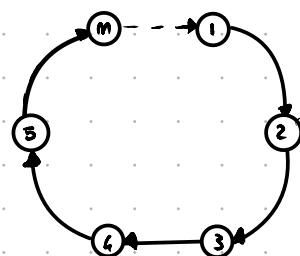
WORST CASE



#MSUs of LEADER

$$\begin{aligned} \# \text{MSUs} &= m + (m-1) + (m-2) + \dots + 1 = \\ &= \sum_{i=1}^m i = \frac{m(m+1)}{2} = \Theta(m^2) \end{aligned}$$

BEST CASE



#MSUs of LEADER

$$\# \text{MSUs} = \underbrace{m + 1 + \dots + 1}_m = 2m = \Theta(m)$$

AVG CASE

- The Prob. for a MSU to make one step is $\frac{1}{m}$, so to complete the circle is $\sum_{i=1}^m i \cdot P(\text{MAKING ONE STEP}) \leq \Theta(m \log m)$
 $\Rightarrow \Theta(m \log m)$

ROUND BASED LEADER ELECTION

- A mode become PASSIVE if its value p is less than the one it receive
- A passive mode only forwards messages

FRANKLIN ALGO (ON BIDIRECTED G)

- Processes send msgs in 2 directions
- At each round at least half Procs become PASSIVE
 (Since every 2 near modes one should win) $\Rightarrow \# \text{Rounds } O(\log m)$
- In each round ACTIVE send 2 msgs, PASSIVE forward 2 msgs $\Rightarrow \# \text{MSUs} = 2m \cdot O(\log m) = O(m \log m)$

27/09/21

ELECTION ON FAULTY SYSTEMS

- We assume that links are faulty free.
- We assume a TIMEOUT, (BOUNDED SYNCHRONICITY)
 - It follows that we can DETECT FAULTY PROCESS
- Assume a fully-connected Topology
- If the leader fail, every process will try to replace it.

IDEA

BULLY ALGO [GARCIA-MOLINA]

Detection Ph.

- TIMEOUT reached while waiting LEADER's ping response

Preparation Ph.

- P_i sends "ELECTION" msg to every $P_j | j > i$

Election Ph.

- If no P_j reply in TIME then I'M THE LEADER \Rightarrow Proclamation Ph

- If $P_j | j > i$ receive P_i msg $\Rightarrow P_j$ reply "YOU CANNOT BE THE LEADER" $\Rightarrow P_j$ (if NOT yet) goes to PREPARATION Ph. and P_i goes to PROCLAMATION Ph. (will wait a LEADER)

Proclamation Ph.

- LEADER P_i LEADER send its "ID" to all $P_j | j < i$

- Other P_j waits for LEADER PROCLAMATION MSG, if TIMEOUT expires \Rightarrow goto PROCLAMATION Ph.

LEADER FAILED DURING
IT'S OWN PROCLAMATION

COMPLEXITY

- When P_i detects LEADER FAILURE it broadcast to his $(m-1)$ neighborhood, getting back at most $(m-1)$ replies

- WORST CASE:

$$\# \text{MSGS}_{\text{worst}} = O(m^3)$$

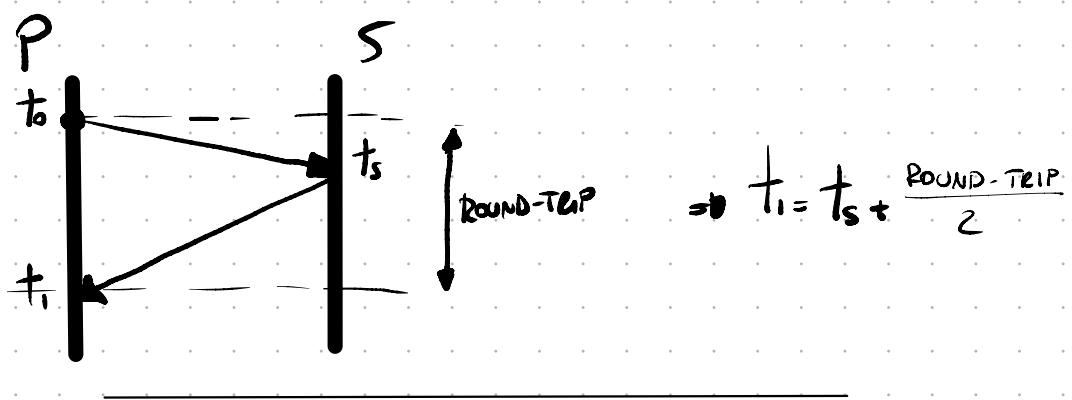
NOTE

- If a failed processor is restarted it tries to elect himself, so if it's the greatest ID proc., we will have 2 LEADERS, we should so consider Epochs to keep the latest elected one.

27/09/21

PHYSICAL CLOCKS

CHRISTIAN'S
PHYSICAL CLOCK



COMPUTER
CLOCK

- Compute the average of the times of every processor; then each one set the average as it's TIME
- Same as COMPUTER CLOCK, but it consider ROUND-TRIP DECAYS too

LOGICAL CLOCKS

LAMPORT'S
LOGICAL
CLOCK

- Each P_i starts with a var. $L = \emptyset$
- BEFORE a LOCAL EVENT $L = L + 1$
- A msg m sent have L as TIMESTAMP RECEIVING m IS AN EVENT
- When a msg m is received $L = \max(L, \text{TIMESTAMP}(m)) + 1$

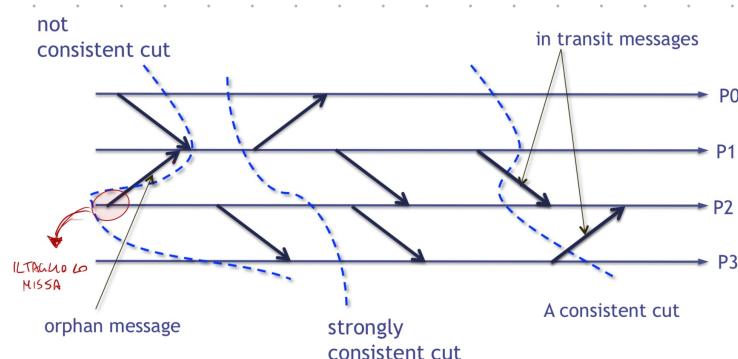
VECTOR
CLOCKS

- P_i start with a vector $V = [V_1, \dots, V_n]$ initialized to \emptyset
 - $V[i]$ contains #events timestamped by P_i
 - $V[j] \neq i$ contains #events of P_j which have interacted with P_i
- BEFORE P_i LOCAL event $V[i] = V[i] + 1$
- A msg sent has the whole vector V
- When a msg m is received from P_j $V[j] = \max(V[j], V_{\text{rcv}}[j]) + 1$

6/10/21

FAULTS AND RECOVERY

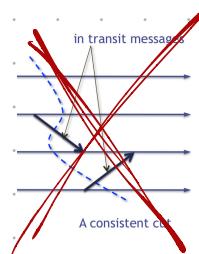
CUT



CHECKPOINT

- Restart from a **CONSISTENT CUT**; So in the cut we save the state of the system.

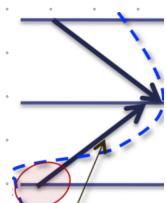
- NB
- In case of RELIABLE system we must use only **STRONGLY CONNECTED CUTS**



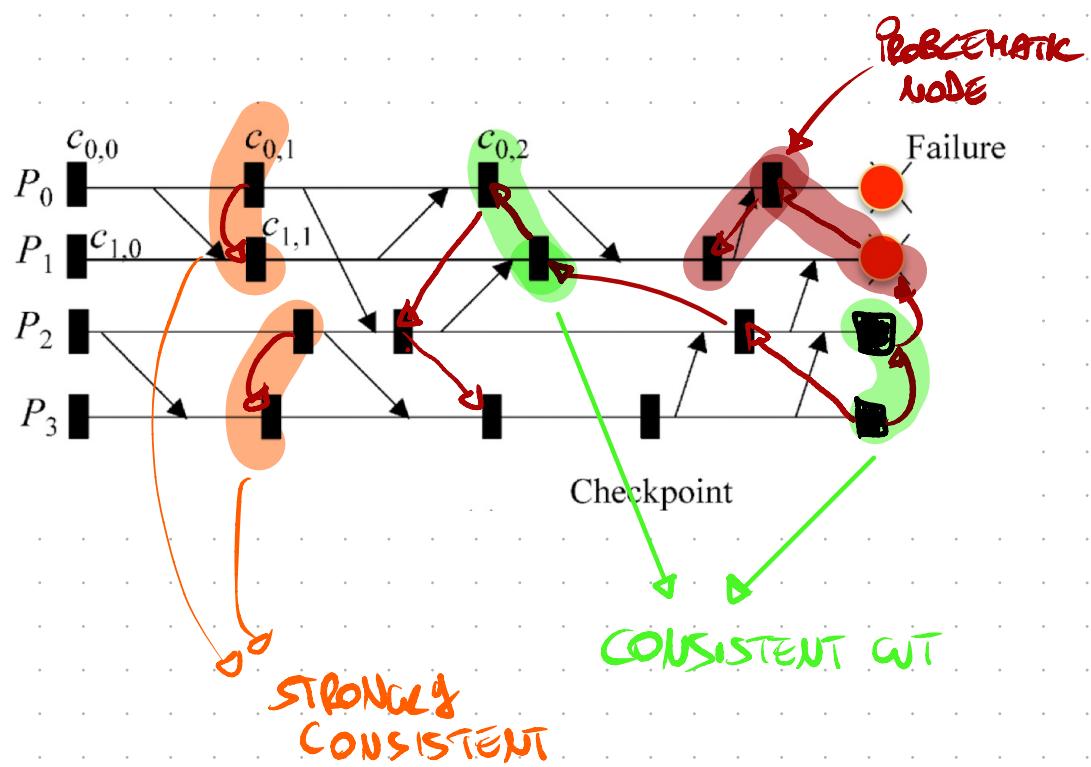
- In this case is not guaranteed that after a restart will be sent the very same msgs.

ORPHAN MSGS

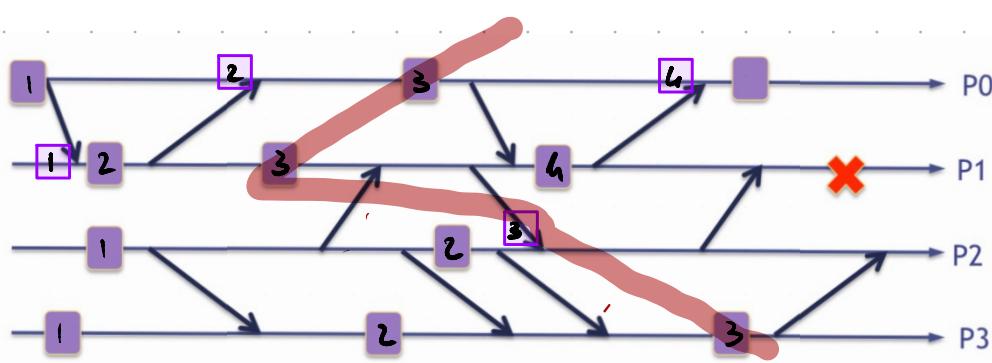
- If it restart the msg will be RECEIVED 2 times from the top processor



ROLLBACK DEPENDENCY GRAPH



CIC



GROUP COMMUNICATION

TYPES OF EVENTS

- 1) REQUEST
- 2) INDICATION (RESPONSE)
- 3) CONFIRMATION (ACK)

```
upon event { Event1 | att11, att21, ... } do
    something
    trigger { Event2 | att12, att22, ... }; // send some event
```

trigger <send | dest, [data1, data2, ...]>

upon event <send | src, [data1, data2, ...]> do

TYPES OF NODE FAILURE

- 1) CRASH STOP
- 2) OMISSION
- 3) CRASH RECOVERY
- 4) BYZANTINE

TYPES OF CHANNEL FAILURE

- 1) FAIR LOSS LINK Non-zero Prob. of delivery
- 2) STUBBORN LINK MSGS DELIVERED INFINITE TIMES
- 3) PERFECT LINK MSGS DELIVERED EXACTLY ONCE

- Request: <flp2pSend | dest, m>
 - Request transmission of message m to node dest
- Indication: <flp2pDeliver | src, m>
 - Deliver message m sent by node src

PROPERTIES

- **FLL** FAIR LOSS
IF m SENT NO w/ NO CRASH \Rightarrow m INFINITELY DELIVERED
- **FCL** FINITE DELIVERY
IF m SENT CO₀ \Rightarrow m DELIVERED CO₀
- **FCL** m NOT DELIVERED IF NOT SENT
No creation

- Request: <sp2pSend | dest, m>
 - Request the transmission of message m to node dest
- Indication: <sp2pDeliver | src, m>
 - deliver message m sent by node src

PROPERTIES

- **SLL** STUBBORN DELIVER
IF m sent to P₁ CORRECT, AND NO CRASH \Rightarrow m DELIVERED CO₀
- **SLL** m DELIVERED \nrightarrow m SENT
No creation

- Implementation
 - Use the Lossy link
 - Sender stores every message it sends in sent
 - It periodically resends all messages in sent

- Request: <pp2pSend | dest, m>
 - Request the transmission of message m to node dest
- Indication: <pp2pDeliver | src, m>
 - deliver message m sent by node src

PROPERTIES

- **PPL** RECEIABLE DELIVERY
IF NO CRASH \Rightarrow m EVENTUALLY DELIVERED
- **PCL** NO DUPLICATION
- **PCL** m DELIVERED AT MOST ONCE
- **PCL** m DELIVERED \nrightarrow m SENT
No creation

- Implementation
 - Use Stubborn links
 - Receiver keeps log of all received messages in Delivered
 - Only deliver (call pp2pDeliver) messages that weren't delivered before

IDEA ALL PERFECTIONLESS SERVER CRASH

Request: <bebBroadcast | m> : Used to broadcast message m to all processes.
Indication: <bebDeliver | src, m> : Used to deliver message m broadcast by process src.

PROPERTIES

- **BEB1** BEST EFFORT VARIETY
IF P₁, P₂ CORRECT \Rightarrow m EVENTUALLY DELIVERED
- **BEB2** NO DUPLICATION
- **BEB3** NO CREATION

- Request: <coBroadcast | m>
- Indication: <coDeliver | src, m>

OPTION \rightarrow CB CASUAL DELIVERY
IF P_i DELIVERS m_i, MUST HAVE DELIVERED ALL MSGS CASUALLY PRECEDING m_i

CASUAL ORDER BROADCAST

TOTALLY DELIVERED BROADCAST

RELIABLE Broadcast

- IF p_i AND p_j BOTH DELIVERS m_1, m_2 , THEY DELIVER THEM IN SAME ORDER.

BEB + RELIABILITY PROPERTY

Request: <rbBroadcast | m> : Used to broadcast message m.

Indication: <rbDeliver | src,m> : Used to deliver message m broadcast by process src.

RB1=Beb1: Validity: If a correct process p, broadcasts a message m, then p, eventually delivers m.

RB2=Beb2: No duplication: No message is delivered more than once.

RB3=Beb3: No creation: If a message m is delivered by some process p_j , then m was previously broadcast by some process p_i .

RB4: Agreement: If a message m is delivered by some correct process p_i , then m is eventually delivered by every correct process p_j .

MUTUAL EXCLUSION

ME1 : Mutual Exclusion

- At most one process can remain in CS at any time
- Safety property

ME2 : Freedom from deadlock

- At least one process is eligible to enter CS
- Liveness property

ME3 : Fairness

- Every process trying to enter must eventually succeed
- Absence of starvation

A measure of fairness: bounded waiting

- Specifies an upper bound on the number of times a process waits for its turn to enter SC \rightarrow n-fairness (n is the MAXIMUM number of rounds)
- FIFO fairness when n=0

LAMPORT ALGO

• "VOGLIO ENTRA"

• "Ack"

• "Ack"

⋮

} TROPPI ACK
MA FUNZIONA

DICART & AGRAWALA

• "VOGLIO ENTRA" +TIMESTAMP

• "Ack" SOLO SE NON VOGLIO ENTRARE

• VOGLIO E HO TIMESTAMP PIÙ ACTO

• USCENDO NANDO "Ack"

TOKEN

TOKEN BING

- NON SERVONO PRESENTATIONI
- APPENA MI ARRIVA IL TOKEN ENTRA



NAIVE ALGO

("IL KING È CHI HA IL TOKEN")

• "VOGLIO ENTRA"

• IL KING METTE IN CODA

• KING ESCE DA CS E PASSA IL TOKEN AL PRIMO IN CODA

⇒ **STARVATION** (POSSE ENTRA DUE VOLTE IN FILA)

⇒ PUOI DIVENTARE KING SENZA AVERLO CHIESTO

SUKIZZI-ZAZZAKI

- OGNI UNO TIENE QUESTA RICHIESTA "VOGLIO ENTRA"
- IL KING HA LA CODA E LA "LAST" (QUESTA MIA GIÀ È STATO KING)
- USCENDO AGGIORNA LA LAST E PASSA TOKEN.

QUORUM BASED

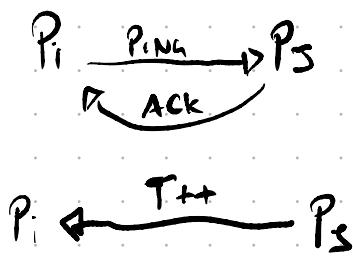
• GRIGLIA

• DEADLOCK

FAILURE DETECTORS

PING ACK

HEARTBEATING



Failure Detectors

("hide/abstract time=delays")

- Basic properties
 - Completeness
 - Every crashed process is eventually suspected by **every** correct process
 - Accuracy
 - No correct process is suspected
- Strong Completeness
 - Every crashed process is eventually suspected by **at least** one correct process
- Weak Accuracy
 - No correct process is **ever** suspected
- Weak Completeness
 - There is **at least** one correct process that is **never** suspected

13

PERFECT FAILURE DETECTOR (ASSUMING SYNCHRONOUS SYSTEM)

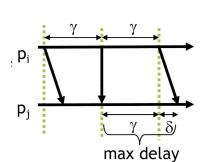
Initialize HBTimeout and DetectTimeout
Upon event (HBTimeout)
For all p_i in P
 Send HeartBeat to p_i
 startTimer (γ , HBTimeout)

Upon event Receive HeartBeat from p_j
 $alive := alive \cup p_j$

Upon event (DetectTimeout)
 $crashed := P \setminus alive$
 for all p_i in $crashed$ Trigger (crashed, p_i)
 $alive := \emptyset$
 startTimer ($\delta + \gamma$, DetectTimeout)

P : set of processes

• MAX DELAY \Rightarrow



EVENTUALLY PERFECT FAILURE DETECTOR (ASYNCHRONOUS SYSTEM)

• ASSUMING THERE IS A FINITE MAX DELAY (PARTIAL & SYNCHRONOUS)

Upon event (HBTimeout)
For all p_i in P
 Send HeartBeat to p_i
 startTimer (γ , HBTimeout)

Upon event Receive HeartBeat from p_j
 $alive := alive \cup p_j$

Upon event (DetectTimeout)
for all p_i in P
 if p_i not in $alive$ and p_i not in suspected
 suspected := suspected $\cup p_i$
 Trigger (suspected, p_i)
 if p_i in $alive$ and p_i in suspected
 suspected := suspected $\setminus p_i$
 Trigger (restore, p_i)
 $T := T + \delta$
 $alive := \emptyset$
 startTimer (T , DetectTimeout)

• When $T > \gamma + \delta$ the system becomes synchronous.

CONSENSUS

PROPERTIES

- Termination
 - Every correct node eventually decides
- Agreement
 - No two correct processes decide differently
- Validity
 - Any value decided is a value proposed
- Integrity:
 - A node decides at most once

- A variant: UNIFORM CONSENSUS
 - Uniform agreement: No two processes decide differently

algorithm I

- A P-based (fail-stop) consensus algorithm
 - The processes exchange and update proposals in rounds and decide on the value of the non-suspected process with the smallest id [Gue95]



- The processes go through rounds incrementally (1 to n); in each round, the process with the id corresponding to that round is the leader of the round
 - The leader of a round decides its current proposal and broadcasts it to all
 - A process that is not leader in a round waits (a) to deliver the proposal of the leader in that round to adopt it, or (b) to suspect the leader

Consensus algorithm II

- A P-based (i.e., fail-stop) uniform consensus algorithm
 - The processes exchange and update proposal in rounds, and after n rounds decide on the current proposal value [Lyn96]



- The “Hierarchical Uniform Consensus” algorithm uses a **perfect failure-detector**, a best-effort broadcast to disseminate the proposal, a **perfect link abstraction** to acknowledge the receipt of a proposal, and a **reliable broadcast** to disseminate the decision
 - Every process maintains a single proposal value that it broadcasts in the round corresponding to its **rank**. When it receives a proposal from a more importantly ranked process (so the hierarchical), it adopts the value
 - In every round of the algorithm, the process whose **rank corresponds to the number of the round** is the leader.
 - A round here consists of two communication steps: within the same round, the leader broadcasts a **PROPOSAL** message to all processes, trying to impose its value, and then expects to obtain an acknowledgment from all correct processes
 - Processes that receive a proposal from the leader of the round adopt this proposal as their own and send an acknowledgment back to the leader of the round
 - If the leader succeeds in collecting an acknowledgment from all processes except detected as crashed, the leader can decide. It disseminates the decided value using a reliable broadcast communication abstraction

GLOBAL STATE COLLECTION

SOME GLOBAL STATES

- { TERMINATION STATE : Useful to start next phase
- DEADLOCK DETECTION
- #MSGS EXCHANGED
- #PROCS AT GIVEN MOMENT

• TERMINATION DETECTION

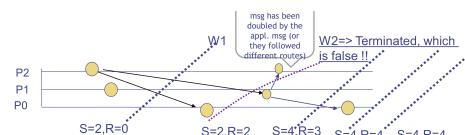
(1) #SENt = #RCVED_Processed

↳ IF THERE ARE ONLY PASSIVE PROCESSES, IT SEEMS TO BE IN TERMINATION STATE (NO MSG EXAG) BUT THOSE PASSIVE CAN BE PROCESSING

- Check if there is a consistent wt. If there is double check, the cut will be waves.

- TO CHECK ACTIVE PASSIVE: If Passive forward the msg then is ACTIVE

- initially each of the n process state is white, Initiator is process 0
 - On msg reception on any Process: state=black;
 - (On each P_i , $i \in [1..n-1]$) On token reception:
 - if state= black token:=1 else token:=token+1
 - state = white; forward token
 - (On P_0) On token reception:
 - if (state=white & token=n) “terminated”, else state=white; forward token=1



- $0 \rightarrow m-1 \rightarrow m-2 \rightarrow \dots \rightarrow 1 \rightarrow 0$
- Ci sono msg che non sono token
- SE MANDI A $ID' > ID_{TUO}$ \Rightarrow DIVENTI NERO
- SE SEI NERO E MANDI IL TOKEN \Rightarrow TOKEN NERO
- SE SEI BIANCO E FORWARDI \Rightarrow RESTI BIANCO
- SE SEI NERO E MANDI MSG A $ID' < ID_{TUO}$ \Rightarrow DIVENTI BIANCO
- SE 'INITIATOR BIANCO' riceve TOKEN BIANCO \Rightarrow TERMINATION

• **RES. DEADLOCK**: When there is a cycle in WAIT-FOR-RES-GRAPH

• **COMMUN. DEAD.**:

Communication deadlock is detected
if the initiator receives ack from all its successors (implying none of its successors can unblock it => it is deadlocked)

DEADLOCK

