

Computer Systems

Exercise Session 11
HS 2024

Approximate Agreement

Approximate Agreement

It enables nodes to obtain values that are:

1. within the range of correct inputs (**correct-range validity**)
 2. ϵ -close for some predefined $\epsilon > 0$ (**ϵ -agreement**)
-
1. $n > 3f$ must hold
 2. synchronous algorithm for $f < n/3$ byzantine nodes
 3. asynchronous algorithm for $f < n/3$ byzantine nodes



Synchronous Approximate Agreement

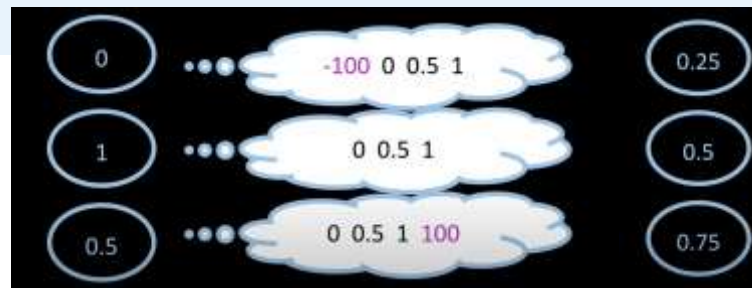
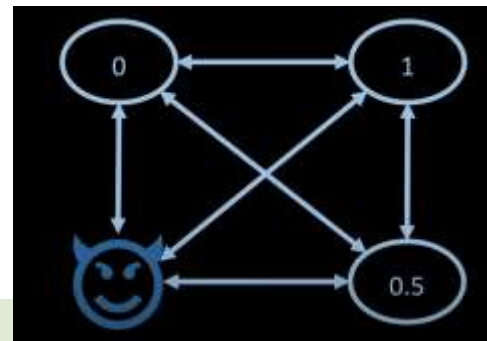
I = a sufficient number of iterations

x_0 = initial value

for $i = 1 \dots I$:

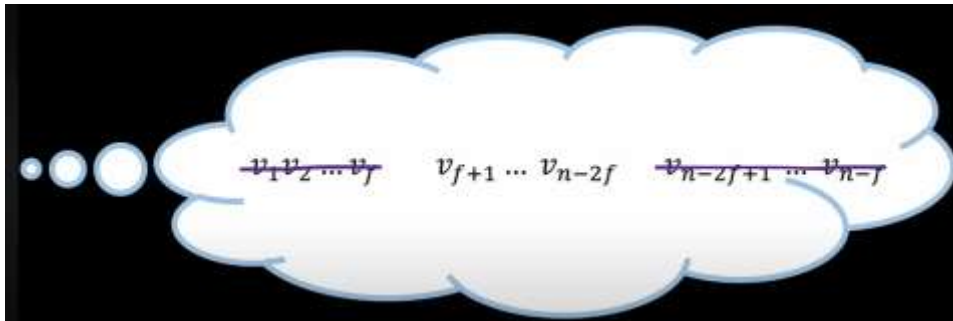
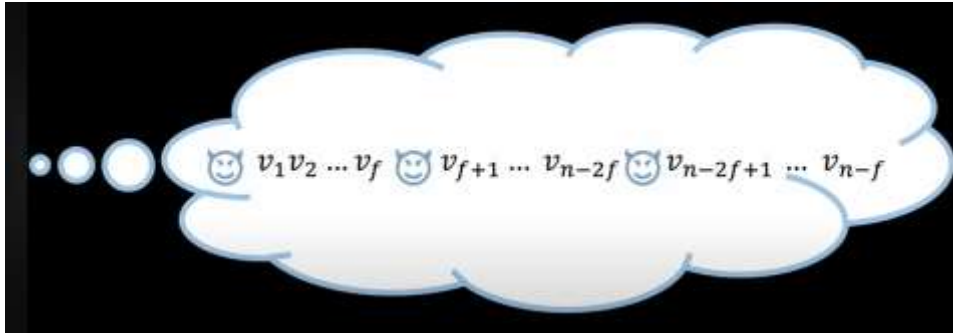
- Distribute your value x_{i-1} .
- R = multiset containing the values received.
- T = multiset containing all but the lowest f and the highest f values in R .
- $x_i = (\min T + \max T) / 2$

Output x_i



Insights

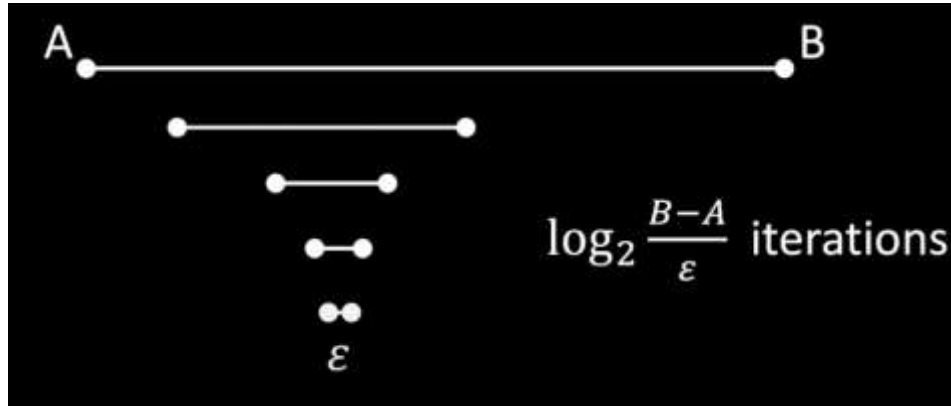
1. The multisets \mathbf{R} contain at most f corrupted values
=> the multisets \mathbf{T} are included in the range of correct values.



Insights

2. If any two correct nodes obtain multisets \mathbf{R} that intersect in $n - f$ values, the range of correct values is **halved** in each iteration

1. **synchronous** model: simply sending your value to everyone is enough.
2. **asynchronous** model: *witness technique*.



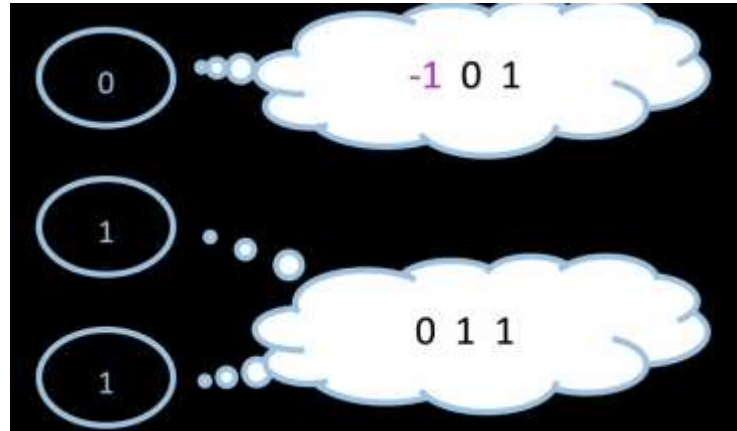
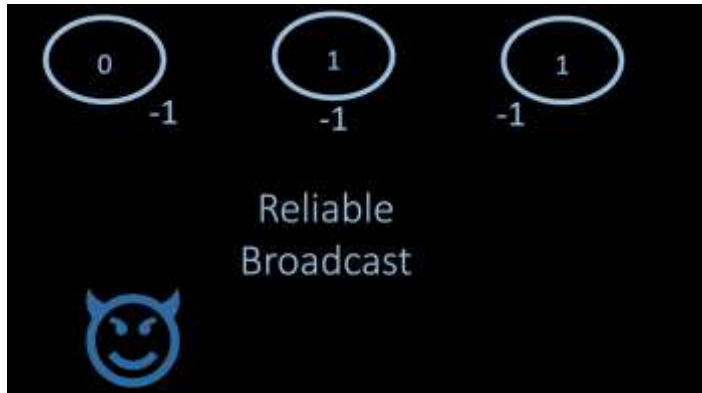
Reliable-Broadcast



- **asynchronous** network with $f < n/3$ byzantine nodes
- Properties:
 - If the sender is correct, all correct nodes accept its value eventually.
 - **If a correct node accepts x , no correct node accepts $y \neq x$.**
 - If a correct node accepts x , all correct nodes accept x eventually.

It's not sufficient with only Reliable-Broadcast

- As nodes only accept n-f values, algorithm may fail due to bad scheduling



Witness Technique



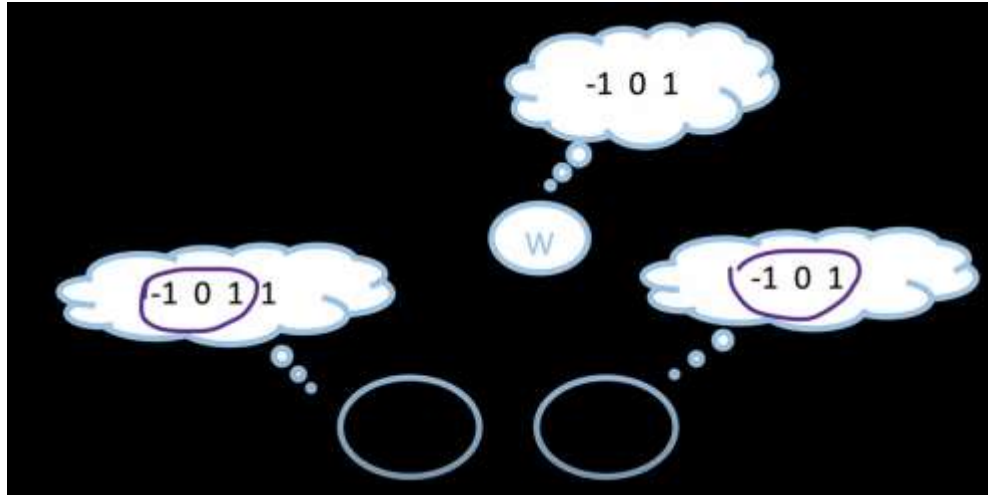
Key idea:

Once a node accepts values from $n - f$ nodes via Reliable Broadcast, it tries to convince all nodes to wait ***a bit longer***: so that they receive these nodes' values as well.

=> nodes obtain multisets \mathbf{R} that pair-wise intersect in $n - f$ values.

Insights

1. Value from a correct witness will be accepted by all correct nodes.
2. Value of correct nodes will be collected in R by all correct nodes.



Algorithm 21.12 The Witness Technique: Iteration i

- 1: Code for node v with input x .
 - 2: Let $R = \emptyset$, $S = \emptyset$, $W = \emptyset$.
 - 3: Send x to all the nodes via Algorithm 18.11 (in the instance for iteration i , with sender v).
 - 4: **upon** accepting $\text{msg}_{i,u}(y)$ from u via Algorithm 18.11 (in the instance for iteration i with sender u):
 - 5: Add y to R and u to S .
 - 6: The first time when $|S| \geq n - f$ holds:
 - 7: Send $\text{wait}_i(S)$ to all the nodes.
 - 8: **end upon**
 - 9: **upon** receiving $\text{wait}_i(S')$ from u such that $|S'| \geq n - f$:
 - 10: When $S' \subseteq S$, add u to W .
 - 11: The first time when $|W| \geq n - f$:
 - 12: Output R .
 - 13: **end upon**
-

Consistency & Logical Time



Consistency models

- Linearizable
- Sequentially Consistency
- Quiescent Consistency

Theorem:

Linearizable implies both sequentially and quiescent consistency.



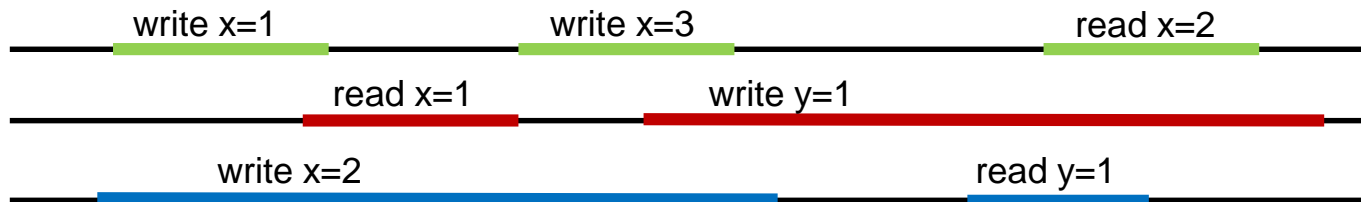
Linearizable

- “one global order”
- Linearizable \rightarrow put points on a “line”
- Strongest assumption, implies other two



Linearizable

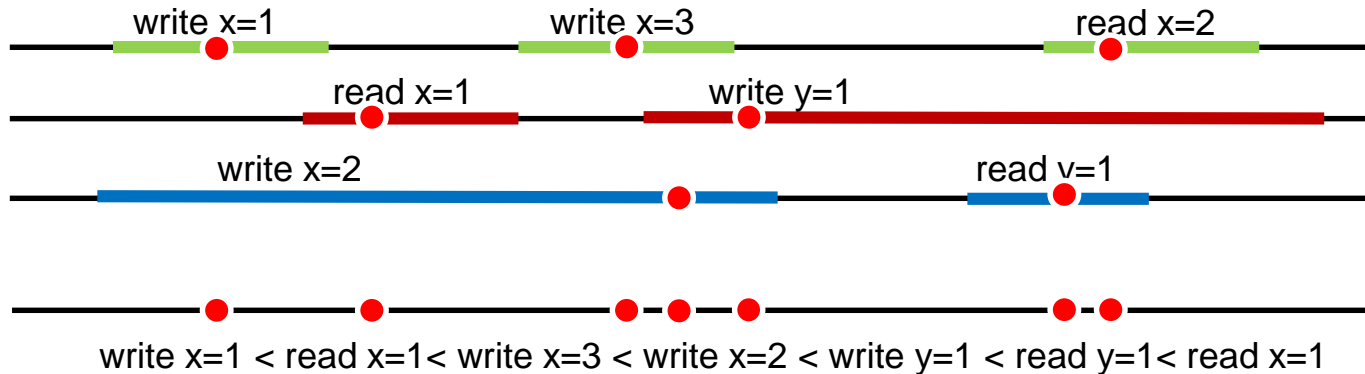
- “one global order”
- Linearizable \rightarrow put points on a “line”
- Strongest assumption, implies other two





Linearizable

- “one global order”
- Linearizable \rightarrow put points on a “line”
- Strongest assumption, implies other two





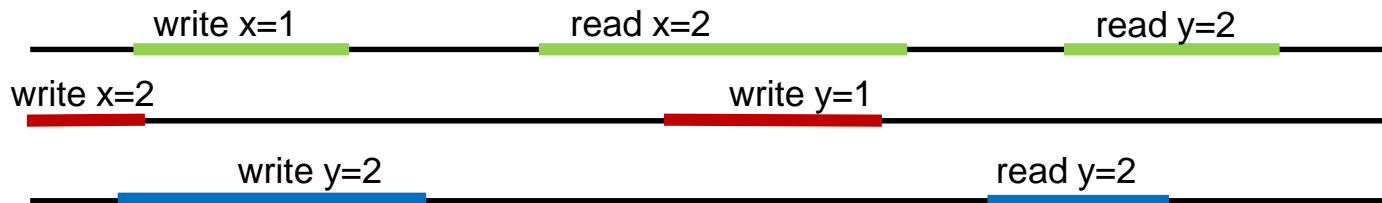
Sequential Consistency

- “per thread order”
- Sequential consistency → build “sequences”



Sequential Consistency

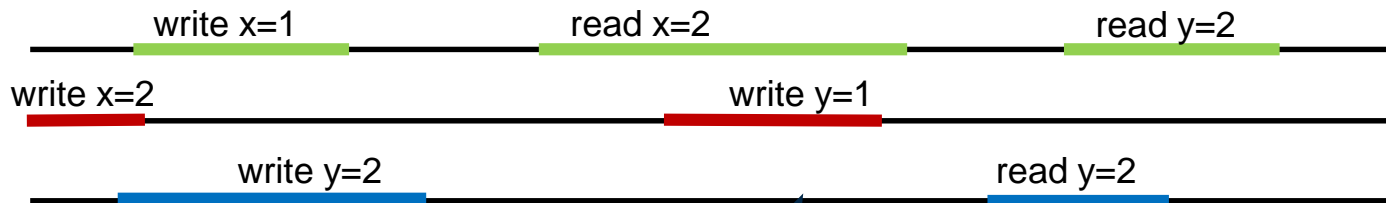
- “per thread order”
- Sequential consistency → build “sequences”





Sequential Consistency

- “per thread order”
- Sequential consistency \rightarrow build “sequences”

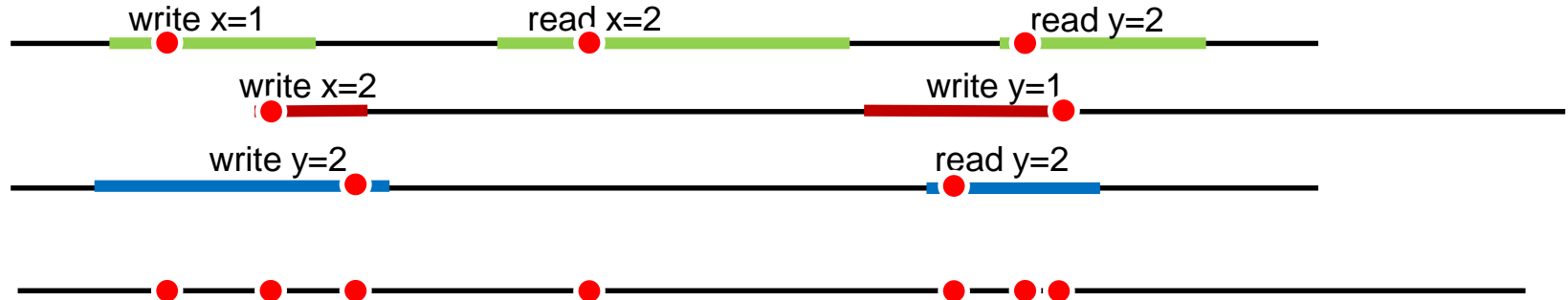


Not
linearizable



Sequential Consistency

- “per thread order”
- Sequential consistency \rightarrow build “sequences”

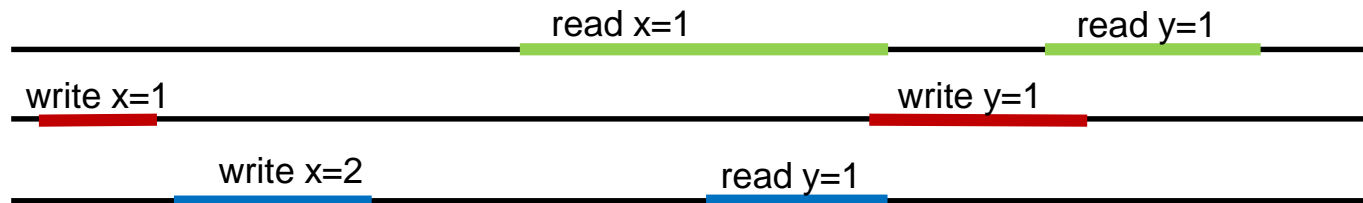


write x=1 < write x=2 < write y=2 < read x=2 < read y=2 < read y=2 < write y=1



Quiescent Consistency

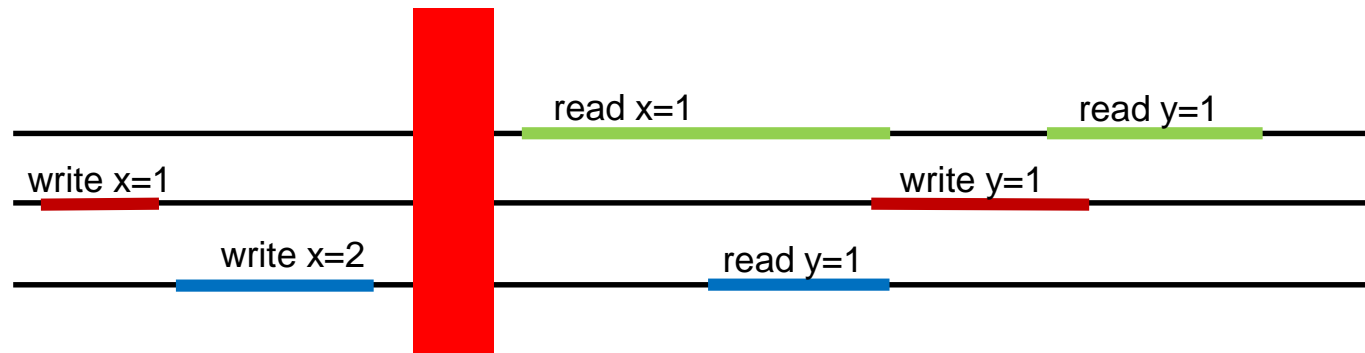
- Synchronizes all threads on quiescent point, i.e. point where no execution happens





Quiescent Consistency

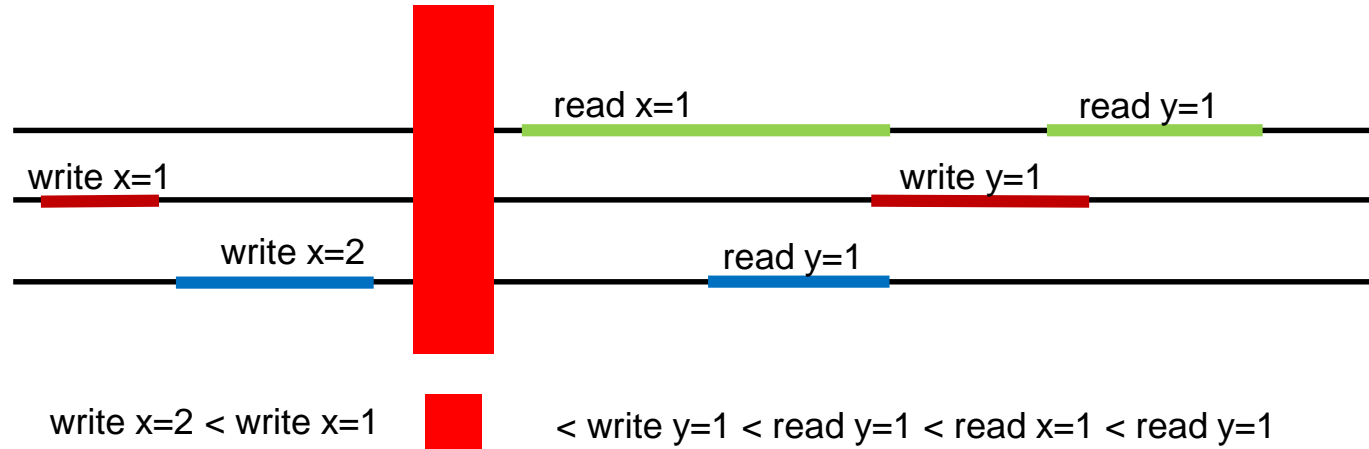
- Synchronizes all threads on quiescent point, i.e. point where no execution happens





Quiescent Consistency

- Synchronizes all threads on quiescent point, i.e. point where no execution happens





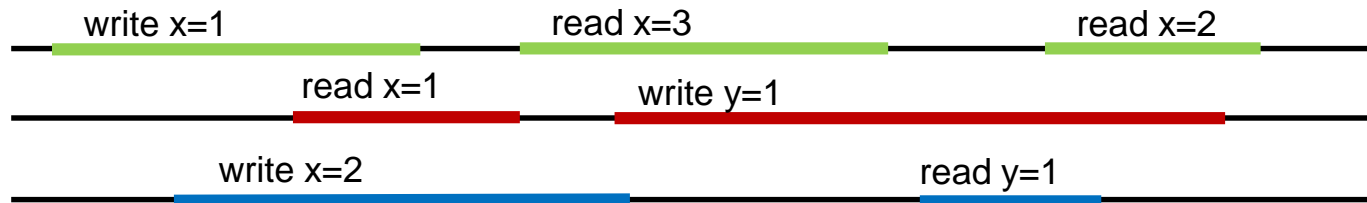
Composable Consistency

- Definition: If you only look at all operations concerning any object and the execution is consistent, then also the whole execution is consistent
- Sequentially consistent is not composable
- Linearizability is composable
- Quiescent consistency is composable



Composable Consistency

- Definition: If you only look at all operations concerning any object and the execution is consistent, then also the whole execution is consistent
- Sequentially consistent is not composable
- Linearizability is composable
- Quiescent consistency is composable





Logical Clocks:

- Happened before relation " \rightarrow " holds
 - 1) IF $f < g$ on the same node
 - 2) Send happens before receive
 - 3) If $f \rightarrow g$ and $g \rightarrow h$, then $f \rightarrow h$ (transitivity)



Logical Clocks:

- Happened before relation " \rightarrow " holds
 - 1) IF $f < g$ on the same node
 - 2) Send happens before receive
 - 3) If $f \rightarrow g$ and $g \rightarrow h$, then $f \rightarrow h$ (transitivity)
- $C(a)$: timestamp of event a

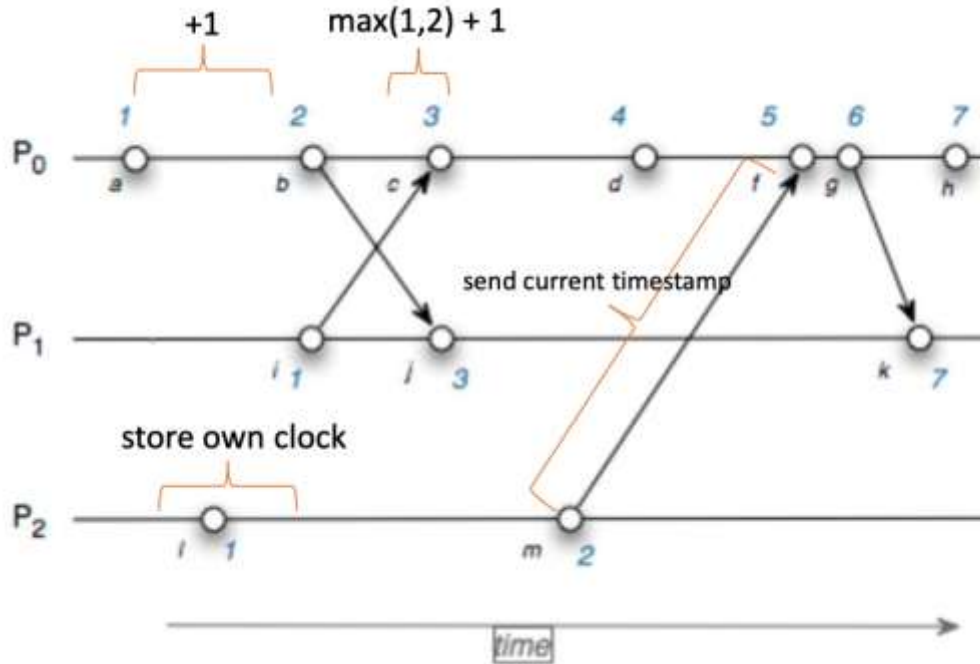


Logical Clocks:

- Happened before relation “ \rightarrow ” holds
 - 1) IF $f < g$ on the same node
 - 2) Send happens before receive
 - 3) If $f \rightarrow g$ and $g \rightarrow h$, then $f \rightarrow h$ (transitivity)
- $C(a)$: timestamp of event a
- **logical clocks: $a \rightarrow b$ implies $c(a) < c(b)$**
- **Strong logical clock: $c(a) < c(b)$ implies $a \rightarrow b$ (in addition)**

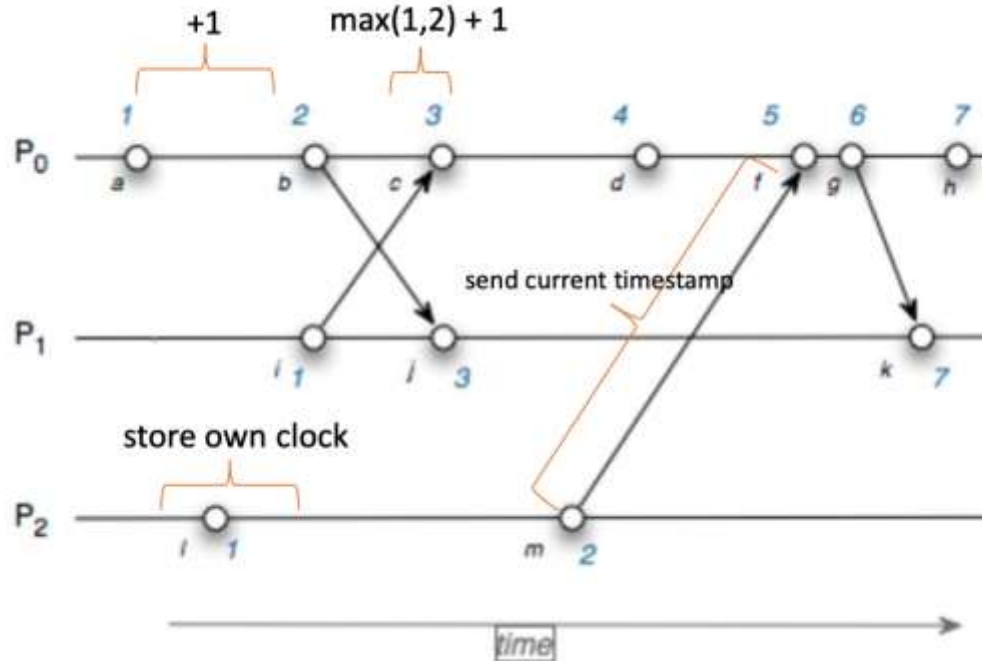


Lamport Clocks:





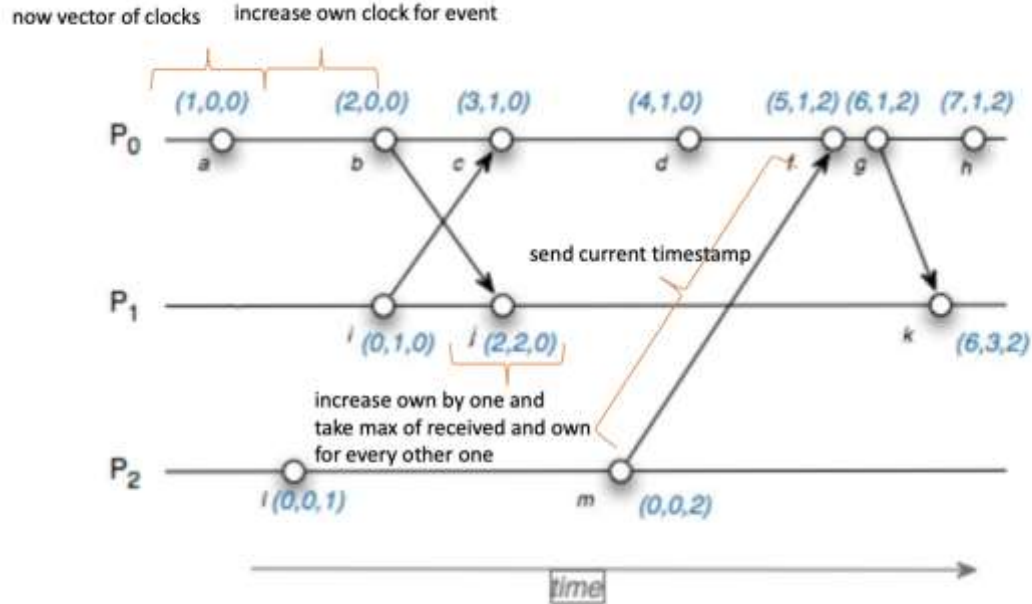
Lamport Clocks:



Weak logical clock: $a \rightarrow b$ implies $c(a) < c(b)$ but not vice versa



Vector Clocks:





Vector Clocks:

- What **does $c(a) < c(b)$ mean now?**
 - if all the entries in $a \leq b$ and at least one entry where $a < b$
- Is a **logical clock** (so if $a \rightarrow b$ then $c(a) < c(b)$)
- Is also a **strong logical clock** (if $c(a) < c(b) \rightarrow a \rightarrow b$)

Intuition: because in order to achieve $c(a) < c(b)$, all entries have to be at least as big, so a message from a must have reached b (not necessarily directly) so that b has the right value



Consistent Snapshot:

- **Cut:** prefix of a distributed execution
- **Consistent Snapshot:**

a cut where for every operation g in that cut, if $f \rightarrow g$, then the cut contains f

→ if all “connected” preceding operations are included

- With number of consistent snapshots, one can make conclusions about degrees of concurrency in system



Quiz

1. Does sequential consistency imply quiescent consistency?
2. Are there guarantees a Lamport clock can achieve a vector clock cannot?
3. Does a high number of consistent snapshots imply a high level of concurrency?



Quiz

1. Does sequential consistency imply quiescent consistency? - **Wrong**

$$x = 2 * x$$

$$x = x + 1$$

e.g. $x=1.5$ is a valid outcome for sequential consistency, but not quiescent

2. Are there guarantees a Lamport clock can achieve a vector clock cannot?

No, because the concept of a Lamport clock is included in the vector clock concept

3. Does a high number of consistent snapshots imply a high level of concurrency? - **True**