



## Where we're at

Last lecture, we discussed classic **distributed algorithms** for commitment and consensus.

Today, we'll discuss algorithms for establishing **global properties** about the system state.

# Global Properties

We'll look at two global properties:

- ① How to tell if a distributed computation has terminated.
- ② How to get a snapshot of the current state of a distributed system.

# Global Properties

We'll look at two global properties:

- ① How to tell if a distributed computation has terminated.
- ② How to get a snapshot of the current state of a distributed system.

## Question

For single-machine systems, these are relatively easy. Why are they hard in a distributed setting?

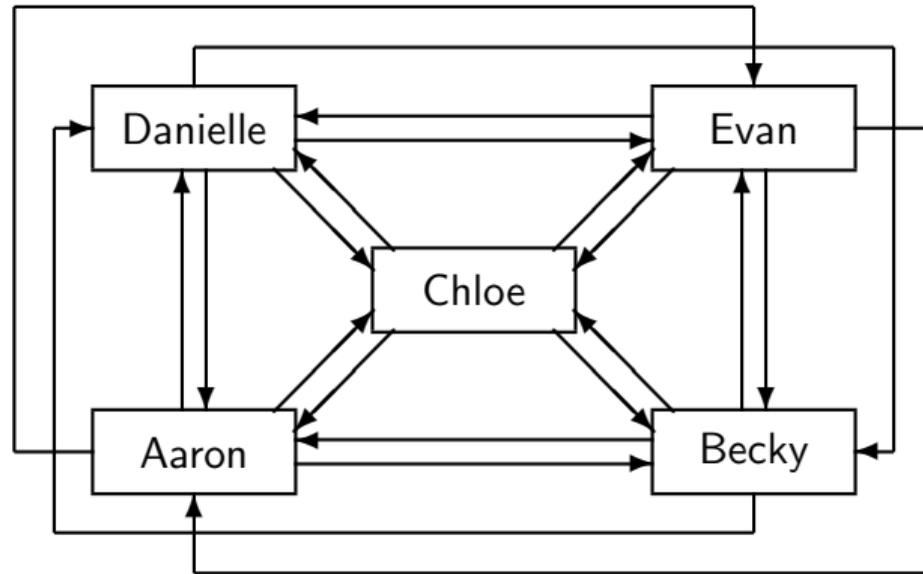
# Setting

## NB

The algorithms we present this lecture are not presented as stand-alone processes.

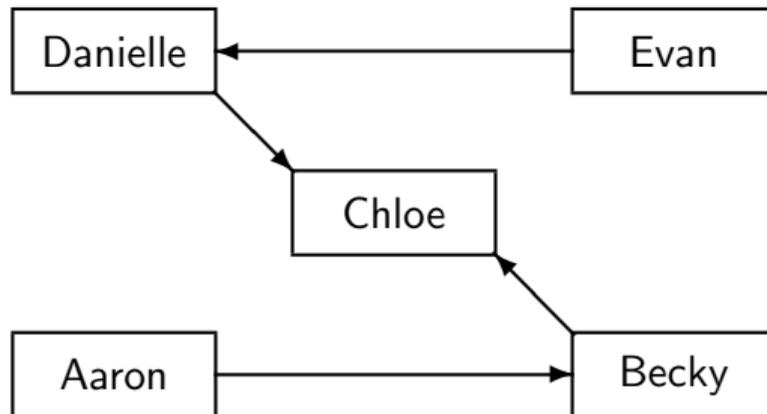
Rather, they are snippets of code meant to be integrated into some *underlying computation* the system is doing. For example, whenever the underlying computation wishes to send a message, the algorithm may require some additional local bookkeeping, and additional structure to the message.

## Setting



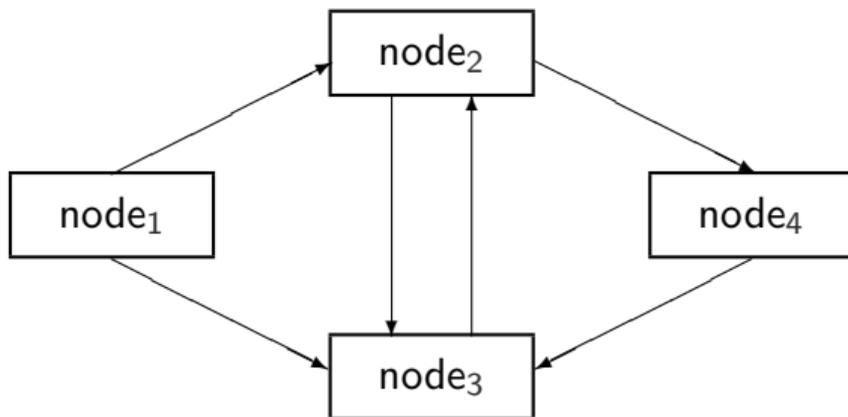
Last week, we assumed **total connectivity**: every node can communicate directly with every other node.

## Setting



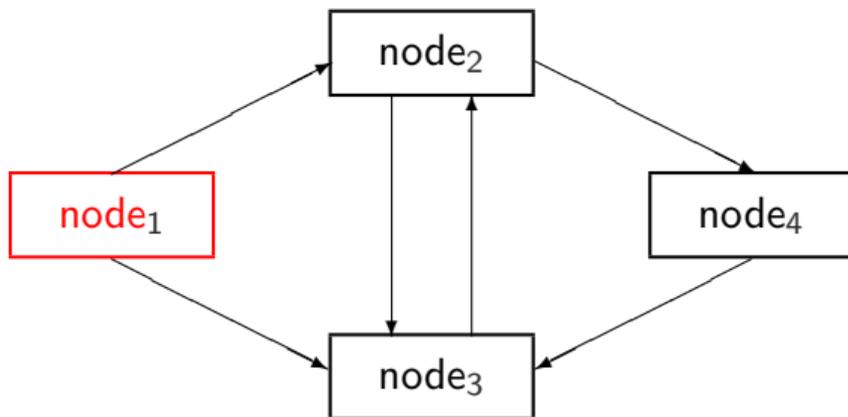
This week, nodes may be partially connected. Connectivity may be asymmetric. Multi-hop communication may be required to reach certain nodes, or there may be no path at all between two particular nodes.

## Distributed System with an Environment Node



We'll assume an **environment node** with no incoming edges. Every other node must be reachable from it (perhaps via multiple hops).

## Distributed System with an Environment Node



We'll assume an **environment node** with no incoming edges. Every other node must be reachable from it (perhaps via multiple hops).

$node_1$  is our environment node (and none of the others even qualify).











**Algorithm 2.1: Dijkstra-Scholten algorithm (preliminary)**

integer array[incoming] inDeficit  $\leftarrow$  [0,...,0]  
integer inDeficit  $\leftarrow$  0, integer outDeficit  $\leftarrow$  0

**send message**

p1: send(message, destination, myID)  
p2: increment outDeficit

**receive message**

p3: receive(message, source)  
p4: increment inDeficit[source] and inDeficit

**send signal**

p5: when inDeficit  $>$  1 or  
    (inDeficit = 1 and isTerminated and outDeficit = 0)  
p6: E  $\leftarrow$  some edge E with inDeficit[E]  $\neq$  0  
p7: send(signal, E, myID)  
p8: decrement inDeficit[E] and inDeficit

**receive signal**

p9: receive(signal, \_)  
p10: decrement outDeficit

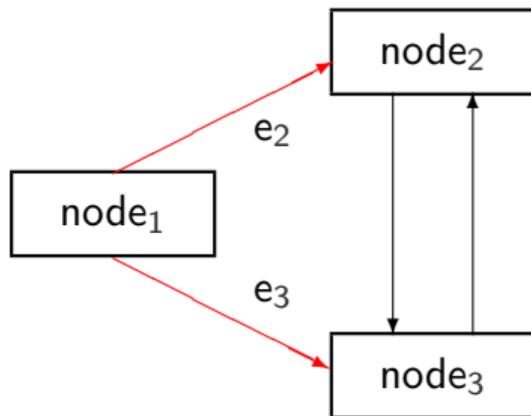
**Algorithm 2.2: Dijkstra-Scholten algorithm (env., preliminary)**integer outDeficit  $\leftarrow$  0**computation**

p1: for all outgoing edges E  
p2:   send(message, E, myID)  
p3:   increment outDeficit  
p4: await outDeficit = 0  
p5: announce system termination

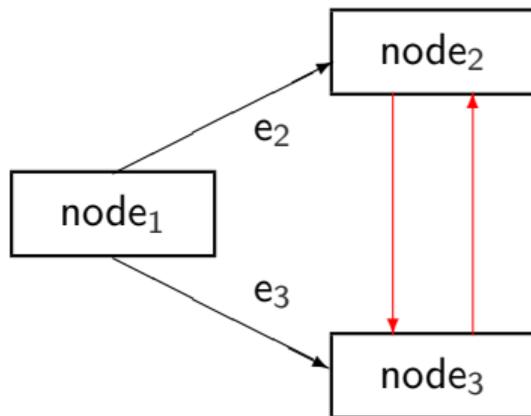
**receive signal**

p6: receive(signal, source)  
p7: decrement outDeficit

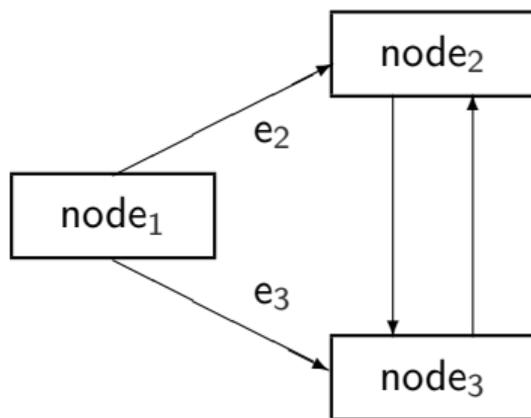
# The Preliminary DS Algorithm is Unsafe



# The Preliminary DS Algorithm is Unsafe



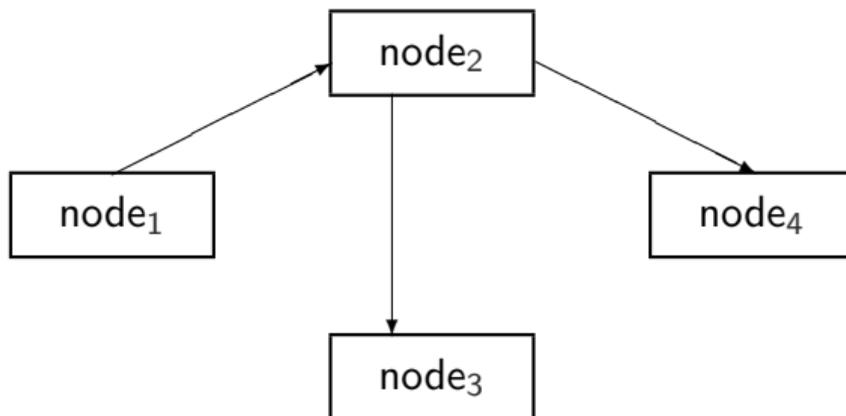
## The Preliminary DS Algorithm is Unsafe



For  $i \in \{2, 3\}$  we have  $\text{node}_i$ 's  $\text{inDeficit} = 2$  and  $\text{inDeficit}[e_i] = 1$  so both can signal back to  $\text{node}_1$ , who is now fooled into announcing termination.

## Spanning Tree

The unsafe example cannot be reconstructed if the channel graph is a tree (with the environment node as root).



**Algorithm 2.3: Dijkstra-Scholten algorithm**

integer array[incoming] inDeficit  $\leftarrow$  [0, ..., 0]

integer inDeficit  $\leftarrow$  0

integer outDeficit  $\leftarrow$  0

integer parent  $\leftarrow$  -1

**send message**

p1: when parent  $\neq$  -1 // Only active nodes send messages

p2: send(message, destination, myID)

p3: increment outDeficit

**receive message**

p4: receive(message, source)

p5: if parent = -1

p6: parent  $\leftarrow$  source

p7: increment inDeficit[source] and inDeficit

**Algorithm 2.3: Dijkstra-Scholten algorithm (continued)****send signal**

p8: when  $\text{inDeficit} > 1$

p9:  $E \leftarrow$  some edge  $E$  for which  
( $\text{inDeficit}[E] > 1$ ) or ( $\text{inDeficit}[E] = 1$  and  $E \neq \text{parent}$ )

p10:  $\text{send}(\text{signal}, E, \text{myID})$

p11: decrement  $\text{inDeficit}[E]$  and  $\text{inDeficit}$

p12: or when  $\text{inDeficit} = 1$  and  $\text{isTerminated}$  and  $\text{outDeficit} = 0$

p13:  $\text{send}(\text{signal}, \text{parent}, \text{myID})$

p14:  $\text{inDeficit}[\text{parent}] \leftarrow 0$

p15:  $\text{inDeficit} \leftarrow 0$

p16:  $\text{parent} \leftarrow -1$

**receive signal**

p17:  $\text{receive}(\text{signal}, \_)$

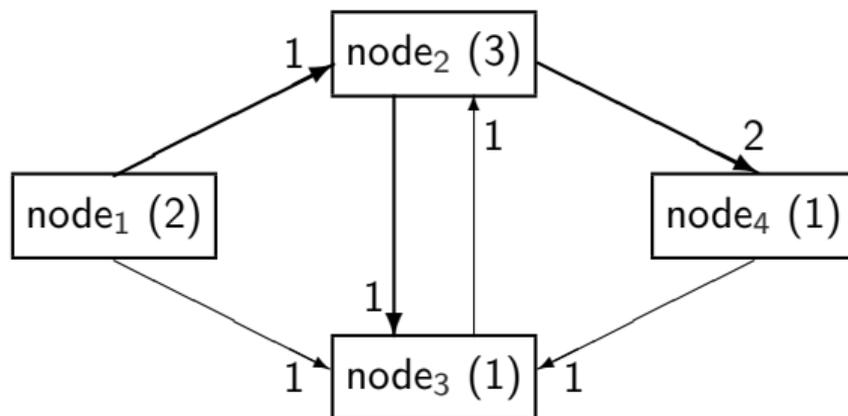
p18: decrement  $\text{outDeficit}$

## Partial Scenario for DS Algorithm

Action	node <sub>1</sub>	node <sub>2</sub>	node <sub>3</sub>	node <sub>4</sub>
1 ⇒ 2	(-1,[ ],0)	(-1,[0,0],0)	(-1,[0,0,0],0)	(-1,[0],0)
2 ⇒ 4	(-1,[ ],1)	(1,[1,0],0)	(-1,[0,0,0],0)	(-1,[0],0)
2 ⇒ 3	(-1,[ ],1)	(1,[1,0],1)	(-1,[0,0,0],0)	(2,[1],0)
2 ⇒ 4	(-1,[ ],1)	(1,[1,0],2)	(2,[0,1,0],0)	(2,[1],0)
1 ⇒ 3	(-1,[ ],1)	(1,[1,0],3)	(2,[0,1,0],0)	(2,[2],0)
3 ⇒ 2	(-1,[ ],2)	(1,[1,0],3)	(2,[1,1,0],0)	(2,[2],0)
4 ⇒ 3	(-1,[ ],2)	(1,[1,1],3)	(2,[1,1,0],1)	(2,[2],0)
	(-1,[ ],2)	(1,[1,1],3)	(2,[1,1,1],1)	(2,[2],1)

$i \Rightarrow k$  means “node<sub>*i*</sub> sends to node<sub>*k*</sub>”; node state notation:  
(parent,inDeficit[E],outDeficit)

## Data Structures After Completion of Partial Scenario



(outDeficit in parentheses, inDeficits on edges)

## Dijkstra-Scholten

Dijkstra-Scholten so far has two shortcomings:

**Traffic overhead** For huge computations, we need to send as many signals as messages

**Unbounded deficits** For huge computations, deficits may grow to the point where they no longer fit in memory.

## Dijkstra-Scholten

Dijkstra-Scholten so far has two shortcomings:

**Traffic overhead** For huge computations, we need to send as many signals as messages

**Unbounded deficits** For huge computations, deficits may grow to the point where they no longer fit in memory.

Unbounded deficit is mainly a theoretical problem. Deficits are only ever stored locally, not sent, so messages are fixed size. Just use arbitrary-precision arithmetic to represent them locally.

## Dijkstra-Scholten

One way to reduce **traffic overhead** in Dijkstra-Scholten is to reduce the deficit as much as possible in a signal. That is, let signals carry a number:

`send(signal, E, myID, N)`

Where  $N$  is the amount deficit we're discharging. (Note the tradeoff: less message quantity, more message complexity)

## Credit-recovery algorithms

In a credit-recovery algorithm, the environment node initially holds 1 credit. Every other node holds 0 credits.

## Credit-recovery algorithms

In a credit-recovery algorithm, the environment node initially holds 1 credit. Every other node holds 0 credits.

The credit is a divisible token. Every active node needs to hold  $> 0$  credit.

## Credit-recovery algorithms

In a credit-recovery algorithm, the environment node initially holds 1 credit. Every other node holds 0 credits.

The credit is a divisible token. Every active node needs to hold  $> 0$  credit.

When you message somebody, give them some credit.

## Credit-recovery algorithms

In a credit-recovery algorithm, the environment node initially holds 1 credit. Every other node holds 0 credits.

The credit is a divisible token. Every active node needs to hold  $> 0$  credit.

When you message somebody, give them some credit.

When a node terminates, it gives all credit back to the environment node.

### Algorithm 2.4: Credit-recovery algorithm (environment node)

float weight  $\leftarrow$  1.0

#### computation

p1: **for** all outgoing edges E  
p2:     weight  $\leftarrow$  weight / 2.0  
p3:     send(message, E, myID, weight)  
p4: **await** weight = 1.0  
p5: announce system termination

#### receive signal

p6: receive(signal, w)  
p7: weight  $\leftarrow$  weight + w

### Algorithm 2.5: Credit-recovery algorithm (non-environment node)

constant integer parent  $\leftarrow 0$  // Environment node  
boolean active  $\leftarrow$  false  
float weight  $\leftarrow 0.0$

#### send message

p1: **if** active  
p2:   weight  $\leftarrow$  weight / 2.0  
p3:   send(message, destination, myID, weight)

#### receive message

p4: receive(message, source, w)  
p5: active  $\leftarrow$  true  
p6: weight  $\leftarrow$  weight + w

#### send signal

p7: **when** terminated  
p8:   send(signal, parent, weight)  
p9:   weight  $\leftarrow 0.0$   
p10:  active  $\leftarrow$  false

## Credit-recovery algorithms

Credit-recovery algorithms reduce message overhead in two ways:

- ① You can discharge all your credit in a single message.
- ② Credit goes straight back to the source. No need for multiple hops up the tree.

## Credit-recovery algorithms

Credit-recovery algorithms reduce message overhead in two ways:

- ① You can discharge all your credit in a single message.
- ② Credit goes straight back to the source. No need for multiple hops up the tree.

### Question

What are some (potential) drawbacks of credit-recovery algorithms?

## Snapshots

A *global snapshot* is a recording of the states of all nodes and channels in the system. This recording is (necessarily) distributed.

Node states record

- values of local variables
- sequences of messages sent and received

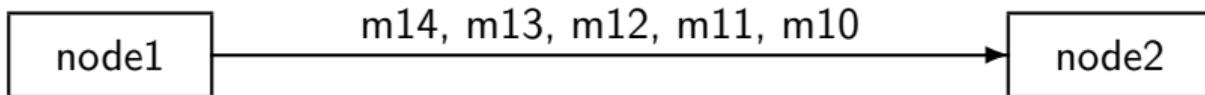
Channel states record

- sequences of messages still in transit

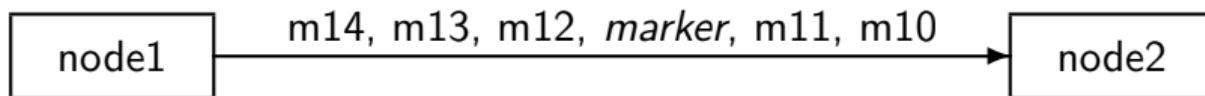
### Definition

A global snapshot is *consistent* iff every sent message is either in transit or received.

## Messages on a Channel



## Sending a Marker



**Algorithm 2.6: Chandy-Lamport algorithm for global snapshots**

integer array[outgoing] lastSent  $\leftarrow [0, \dots, 0]$

integer array[incoming] lastReceived  $\leftarrow [0, \dots, 0]$

integer array[outgoing] stateAtRecord  $\leftarrow [-1, \dots, -1]$

integer array[incoming] messageAtRecord  $\leftarrow [-1, \dots, -1]$

integer array[incoming] messageAtMarker  $\leftarrow [-1, \dots, -1]$

**send message**

p1: send(message, destination, myID)

p2: lastSent[destination]  $\leftarrow$  message

**receive message**

p3: receive(message, source)

p4: lastReceived[source]  $\leftarrow$  message

where message  $\neq$  marker

**NB**

This algorithm assumes all channels are FIFO.

**Algorithm 2.6: Chandy-Lamport algorithm for global snapshots (continued)****receive marker**

```
p6: receive(marker, source)
p7: messageAtMarker[source] ← lastReceived[source]
p8: if stateAtRecord = [-1,...,-1] // Not yet recorded
p9:   stateAtRecord ← lastSent
p10:  messageAtRecord ← lastReceived
p11:  for all outgoing edges E
p12:    send(marker, E, myID)
```

**record state**

```
p13: await markers received on all incoming edges
p14: recordState
```

## The final state

When all marker messages have been received, the final state consists of the following:

- `stateAtRecord[E]`: the last message sent on each outgoing edge E.

## The final state

When all marker messages have been received, the final state consists of the following:

- $\text{stateAtRecord}[E]$ : the last message sent on each outgoing edge  $E$ .
- $\text{messageAtRecord}[E]$ : the last message received on each incoming edge  $E$ .

## The final state

When all marker messages have been received, the final state consists of the following:

- $stateAtRecord[E]$ : the last message sent on each outgoing edge  $E$ .
- $messageAtRecord[E]$ : the last message received on each incoming edge  $E$ .
- The messages in transit on the edge  $E$  are:
  - none, if  $messageAtMarker[E]$  and  $messageAtRecord[E]$  are equal;
  - the messages from  $messageAtRecord[E]+1$  to  $messageAtMarker[E]$ , otherwise.

## The final state

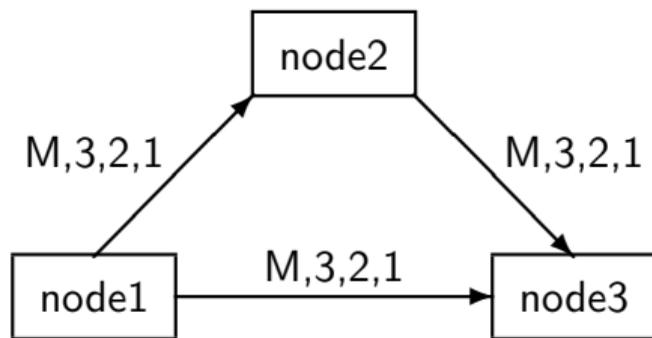
When all marker messages have been received, the final state consists of the following:

- $\text{stateAtRecord}[E]$ : the last message sent on each outgoing edge  $E$ .
- $\text{messageAtRecord}[E]$ : the last message received on each incoming edge  $E$ .
- The messages in transit on the edge  $E$  are:
  - none, if  $\text{messageAtMarker}[E]$  and  $\text{messageAtRecord}[E]$  are equal;
  - the messages from  $\text{messageAtRecord}[E]+1$  to  $\text{messageAtMarker}[E]$ , otherwise.

### NB

Here we've assumed messages are numbered in the order they were sent, to simplify the presentation. Adaptations to record message content are straightforward.

## Messages and Markers for a Scenario



## Scenario for CL Algorithm (1)

Here, the three message from node1 to node2 have been received. Thre three messages from node1 to node3, and from node2 to node3, have all been sent but not received.

Action	node1					node2				
	ls	lr	st	rc	mk	ls	lr	st	rc	mk
	[3,3]					[3]	[3]			
1M $\Rightarrow$ 2	[3,3]		[3,3]			[3]	[3]			
1M $\Rightarrow$ 3	[3,3]		[3,3]			[3]	[3]			
2 $\Leftarrow$ 1M	[3,3]		[3,3]			[3]	[3]			
2M $\Rightarrow$ 3	[3,3]		[3,3]			[3]	[3]	[3]	[3]	[3]

## Scenario for CL Algorithm (2)

Action	node3				
	ls	lr	st	rc	mk
$3 \leftarrow 2$					
$3 \leftarrow 2$		[0,1]			
$3 \leftarrow 2$		[0,2]			
$3 \leftarrow 2M$		[0,3]			
$3 \leftarrow 1$		[0,3]		[0,3]	[0,3]
$3 \leftarrow 1$		[1,3]		[0,3]	[0,3]
$3 \leftarrow 1$		[2,3]		[0,3]	[0,3]
$3 \leftarrow 1M$		[3,3]		[0,3]	[0,3]
		[3,3]		[0,3]	[3,3]

## Chandy-Lamport summary

Once consistent local snapshots are taken, they can be collected through a procedure similar to Dijkstra-Scholten.

The snapshot may be a state that never actually occurred in the system execution. But it's guaranteed to be a state that *could* have occurred in a different interleaving.

## What now?

We're on the home stretch! Next week is recap and exam prep.

In the meantime, **please fill in the myExperience survey** to evaluate the course. Your feedback is tremendously helpful.