# Vector Clocks & Distributed snapshots

CS 452

# Logistics

Problem Set 1 posted: due on Jan 27th

No class on Monday (holiday) and Wednesday (I'm out of town)

# Vector clocks

Precisely represent transitive causal relationships

$T$(A) < $T$(B) <-> *happens-before*(A, B)

Idea: track events known to each node, *on each node*

Used in practice for eventual and causal consistency

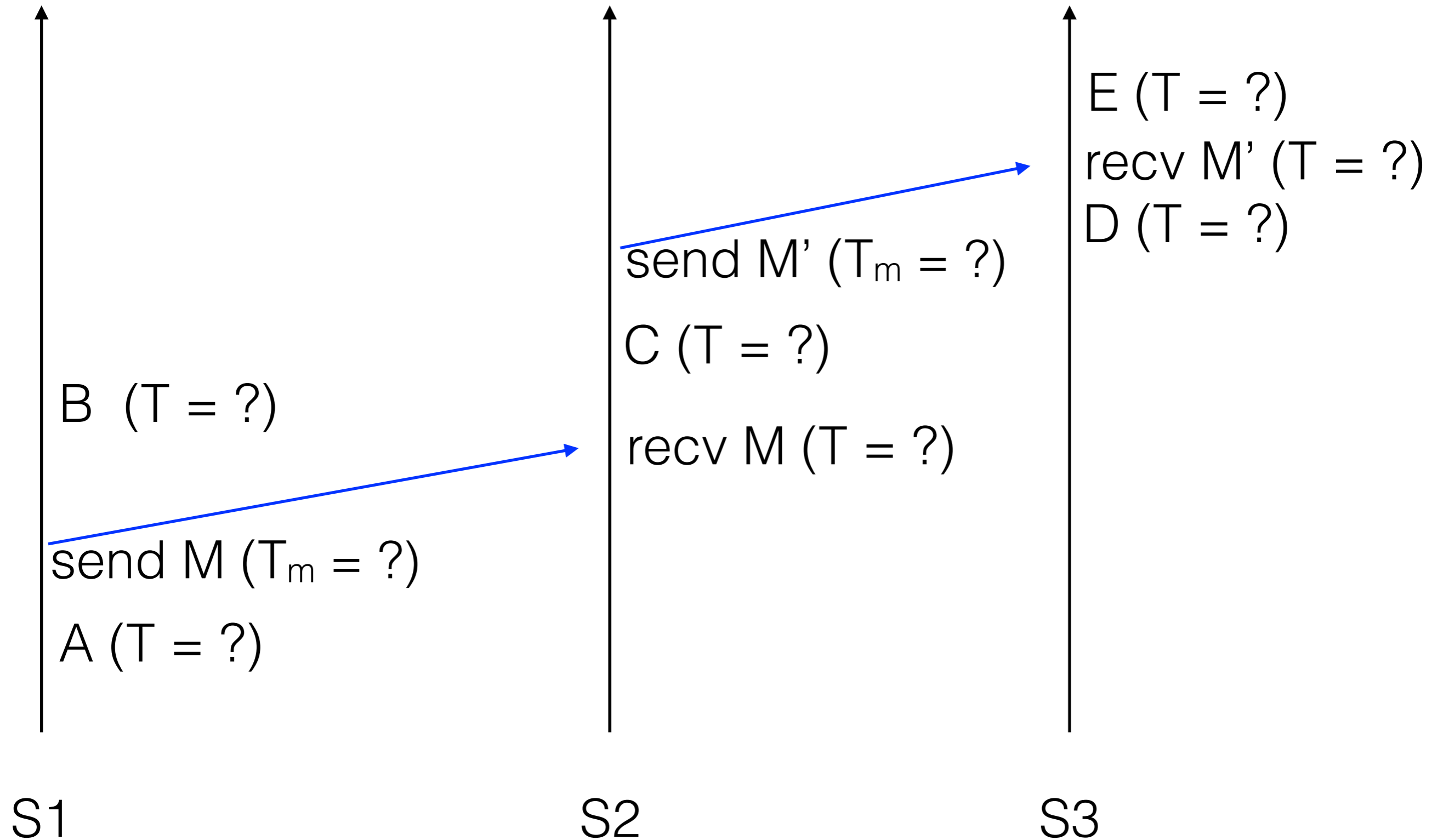- git, Amazon Dynamo, ...

# Vector clocks

Clock is a vector C, length = # of nodes
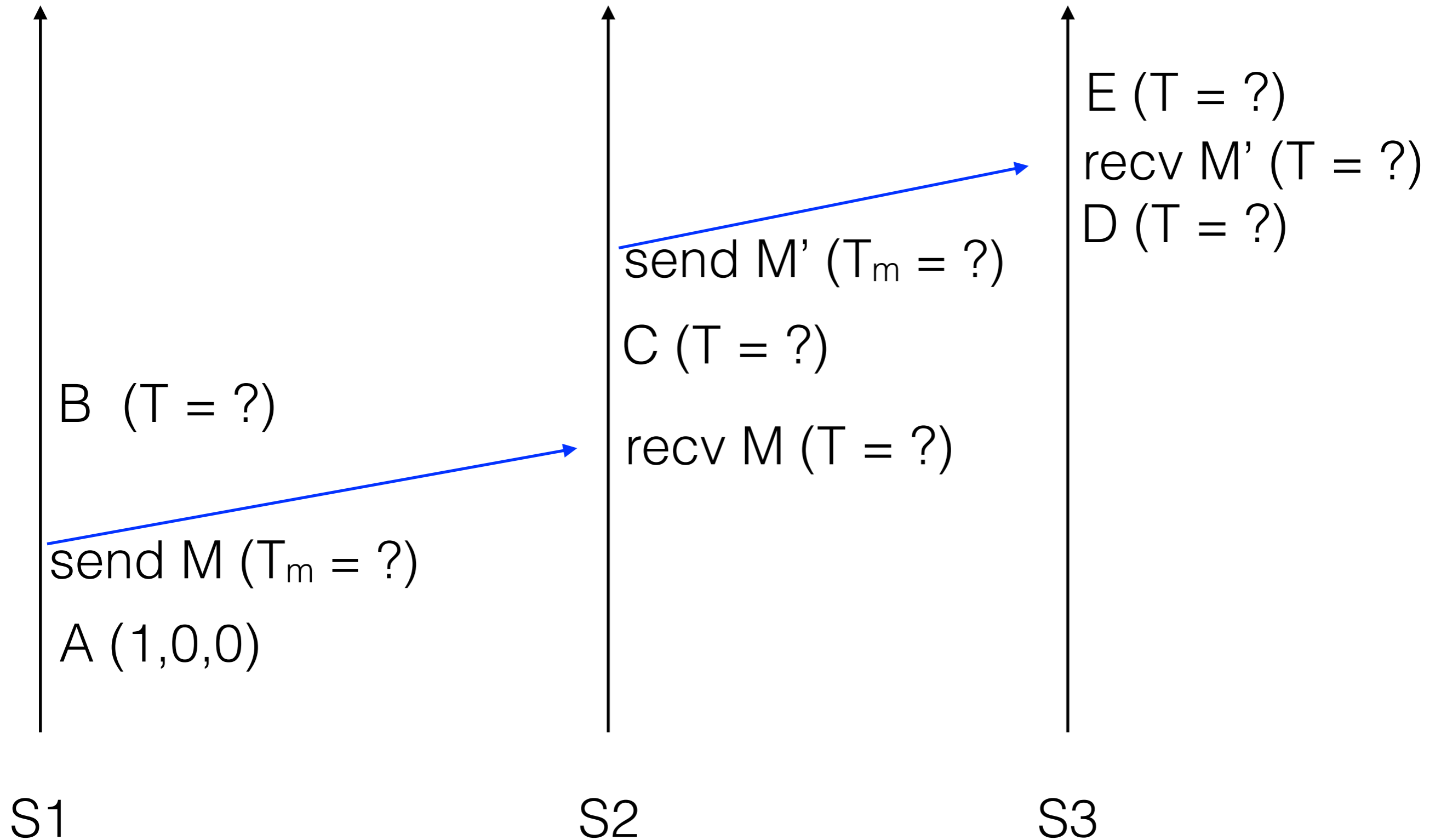
On node i, increment C[i] on each event

On receipt of message with clock $C_m$ on node i:

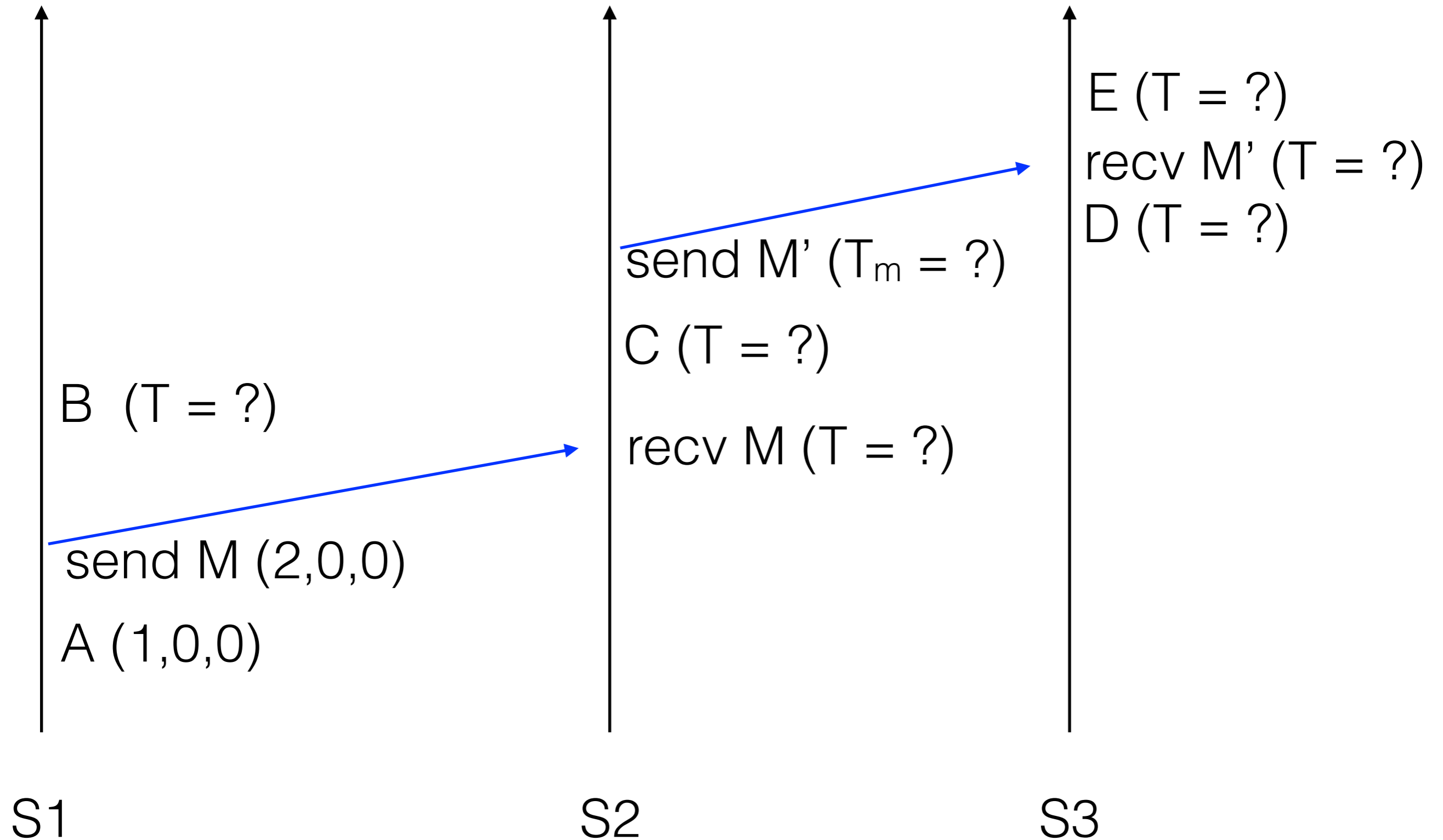- increment C[i]

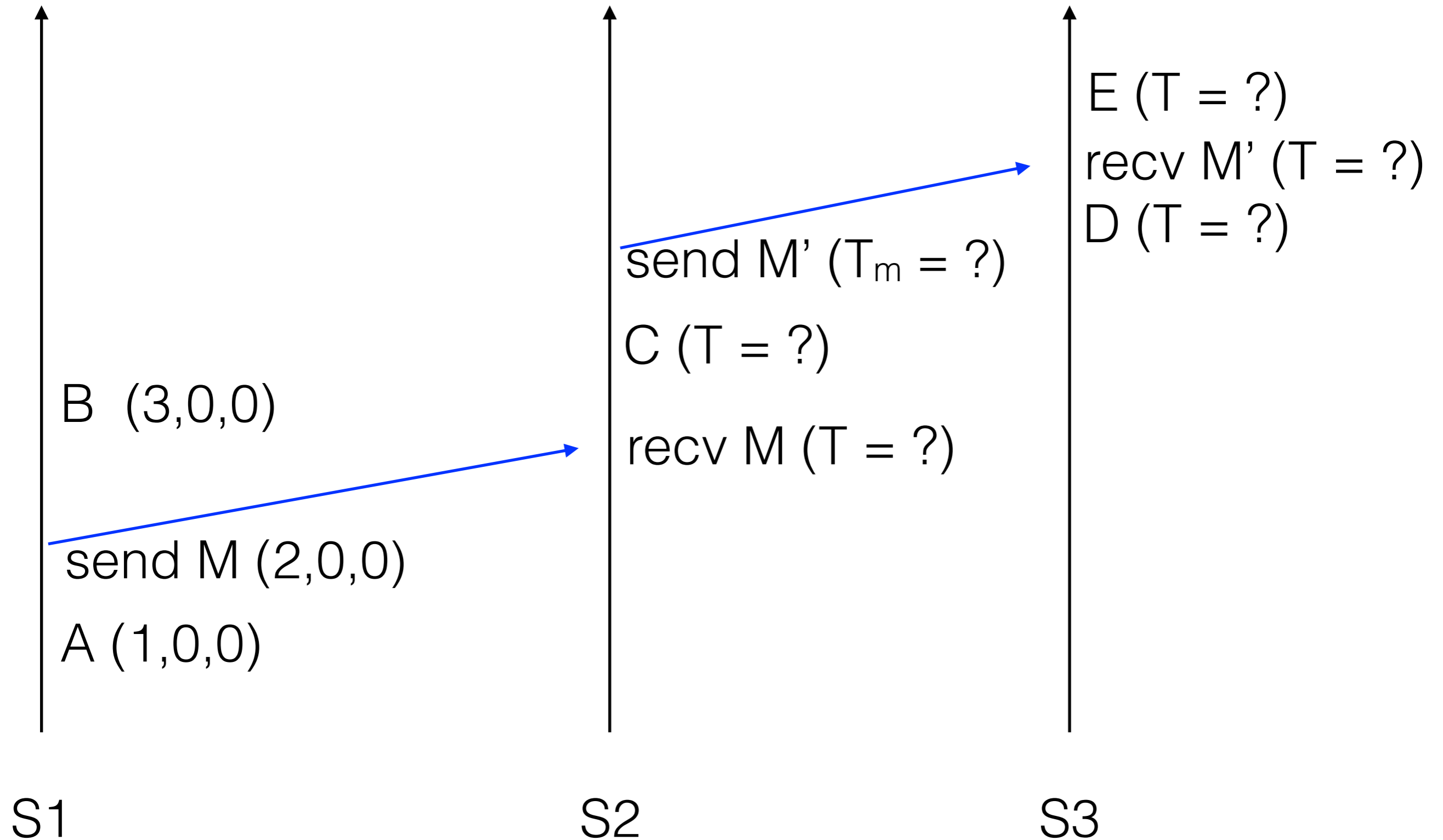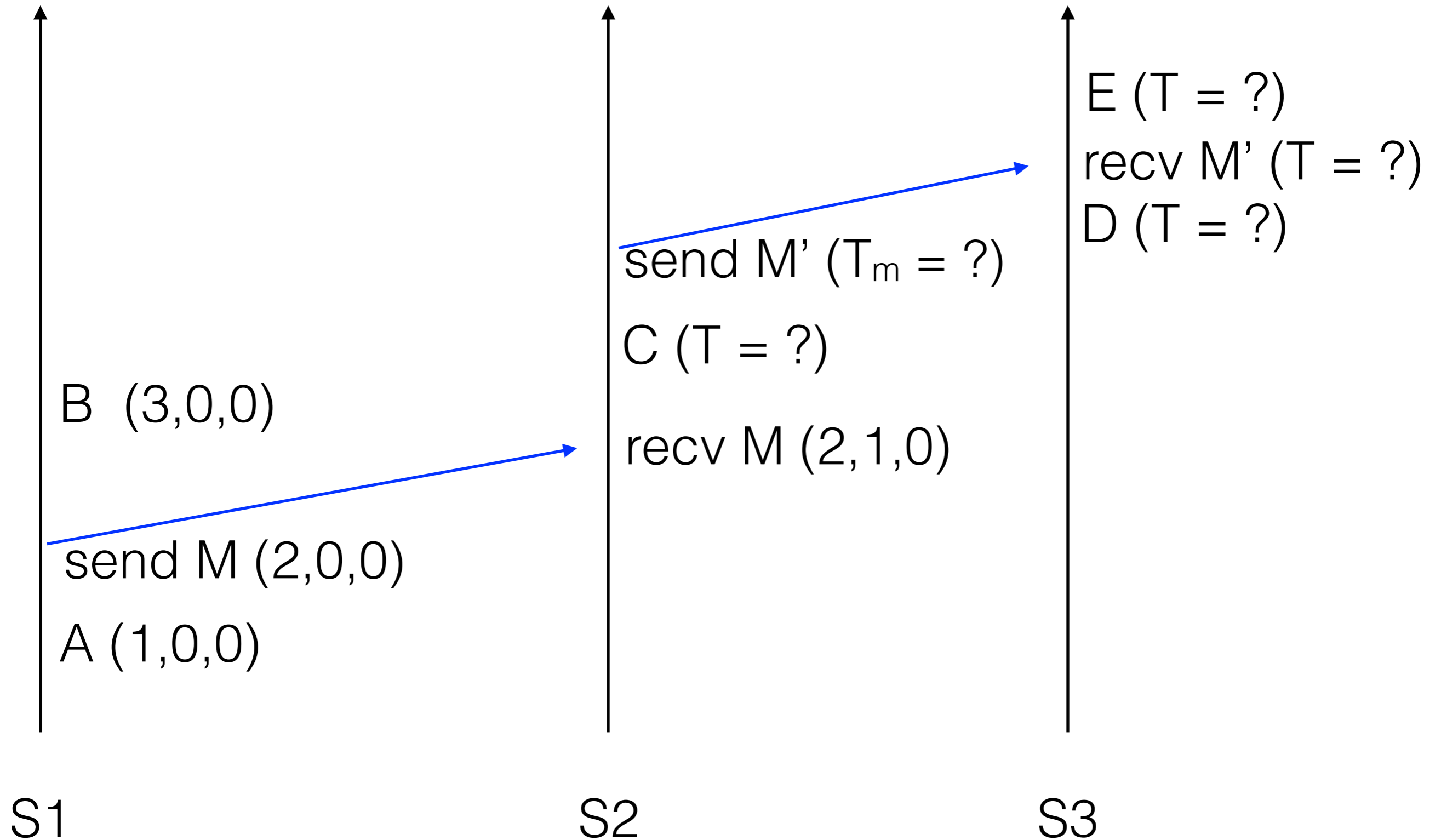- for each j != i

- C[j] = $max$(C[j], $C_m$[j])

# Example



E (T = ?)
recv M' (T = ?)
D (T = ?)

send M' ($T_m$ = ?)

C (T = ?)

recv M (T = ?)

B  (T = ?)

send M ($T_m$ = ?)

A (T = ?)

S1                    S2                    S3

# Example

# Example

E (T = ?)

recv M' (T = ?)

D (T = ?)

send M' ($T_m$ = ?)

C (T = ?)

recv M (T = ?)

B  (T = ?)

send M (2,0,0)

A (1,0,0)

S1                                        S2                                        S3

# Example



S1

B  (3,0,0)

send M (2,0,0)

A (1,0,0)

S2

send M' (T_m = ?)

C (T = ?)

recv M (T = ?)

S3

E (T = ?)

recv M' (T = ?)

D (T = ?)

# Example



E (T = ?)

recv M' (T = ?)

D (T = ?)

send M' ($T_m$ = ?)

C (T = ?)

recv M (2,1,0)

B  (3,0,0)

send M (2,0,0)

A (1,0,0)

S1

S2

S3

# Example

# Example



S1

A (1,0,0)

send M (2,0,0)

B (3,0,0)

S2

recv M (2,1,0)

C (2,2,0)

send M' (2,3,0)

S3

D (T = ?)

recv M' (T = ?)

E (T = ?)

# Example

E (T = ?)

recv M' (T = ?)

D (0,0,1)

send M' (2,3,0)

C (2,2,0)

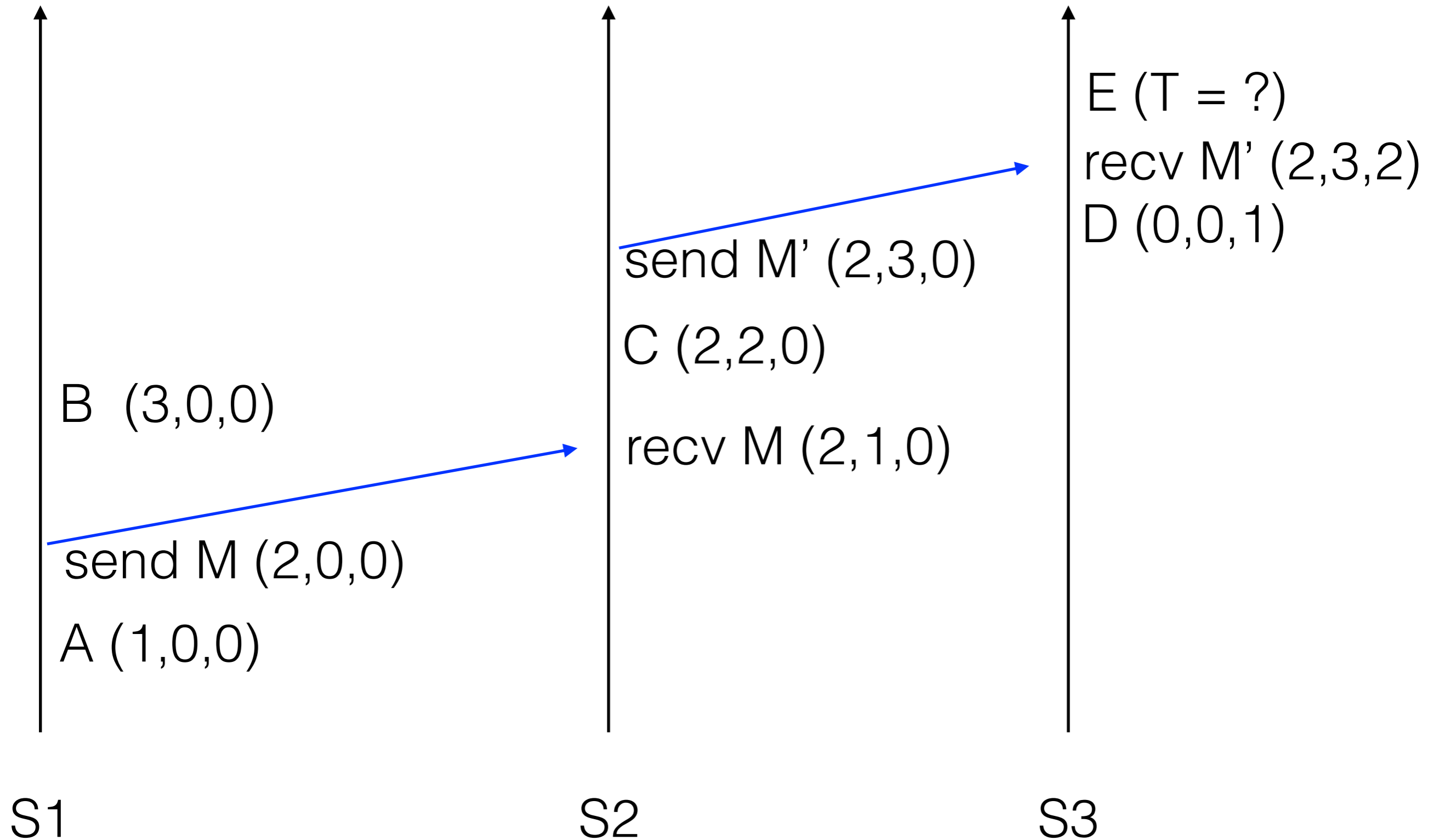recv M (2,1,0)

B  (3,0,0)

send M (2,0,0)
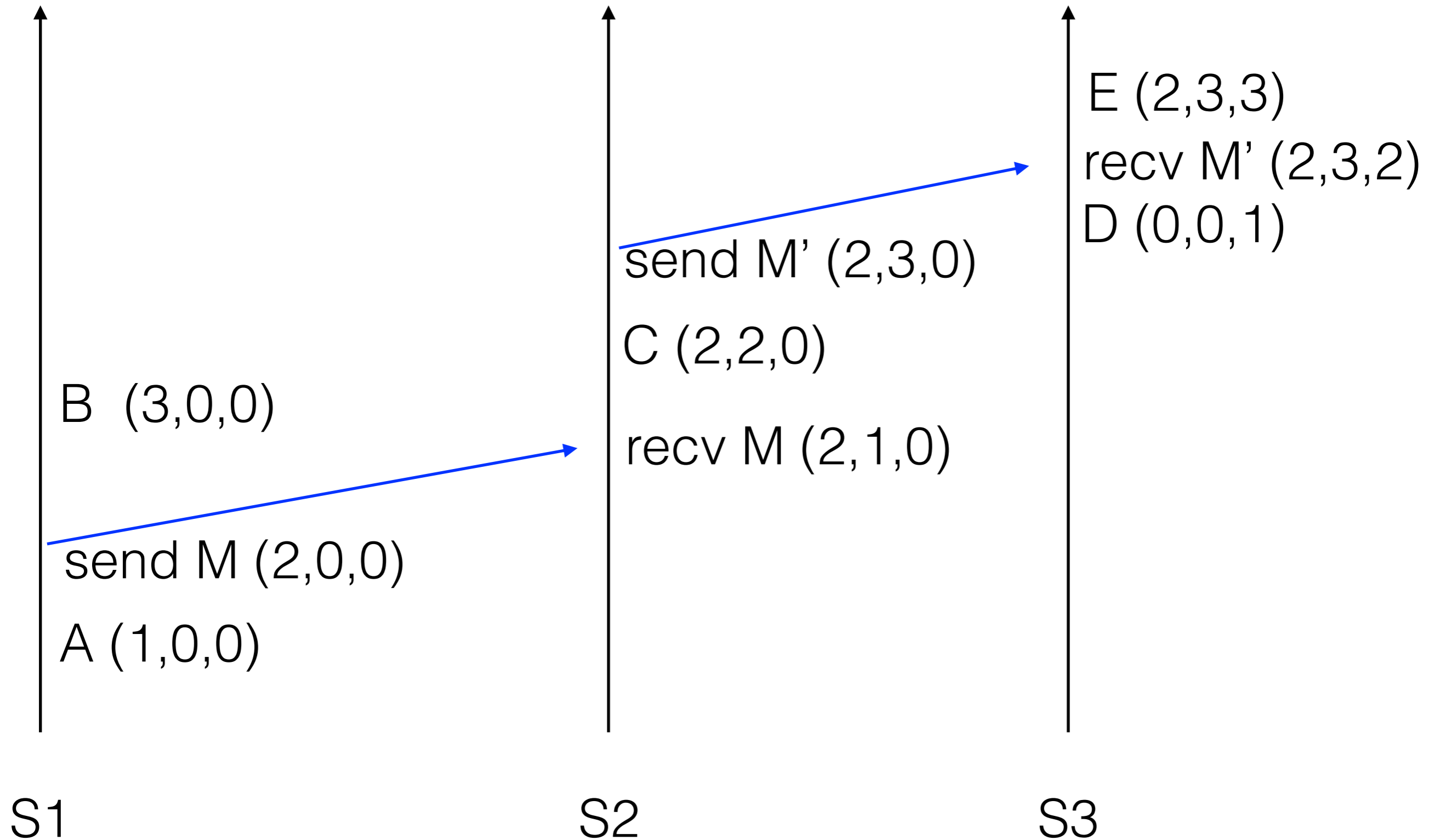
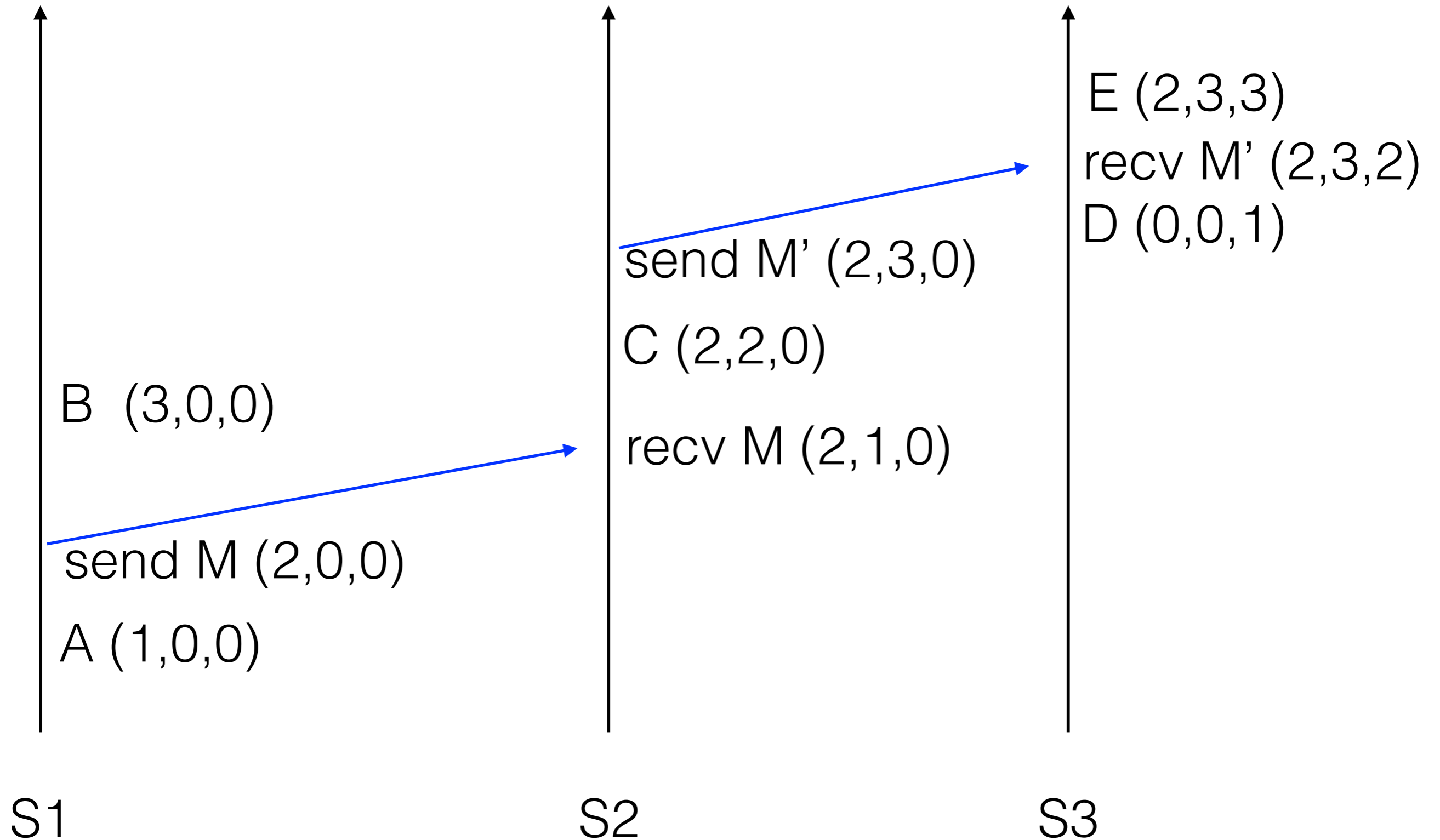A (1,0,0)

S1

S2

S3

# Example

# Example

# Example

# Vector Clocks

Compare vectors element by element
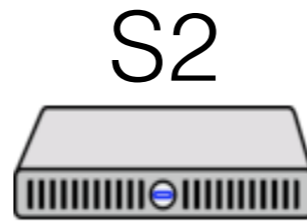
Provided the vectors are not identical,

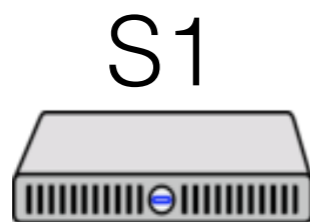If $C_x[i] < C_y[i]$ and $C_x[j] > C_y[j]$ for some i, j

$C_x$ and $C_y$ are concurrent
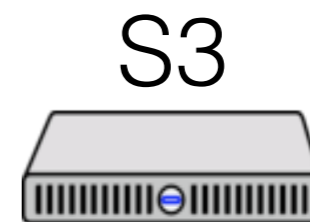

if $C_x[i] <= C_y[i]$ for all i
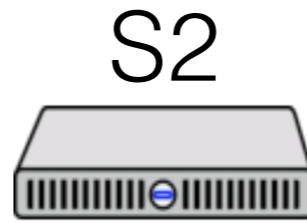
$C_x$ happens before $C_y$

S2

Timestamp: 0
Queue: [S1@0]
$S1_{max}$: 0
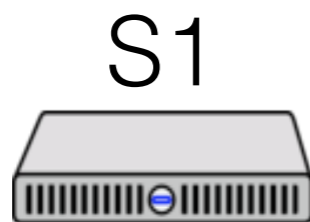$S3_{max}$: 0

S1

Timestamp: 0
Queue: [S1@0]
$S2_{max}$: 0
$S3_{max}$: 0

S3

Timestamp: 0
Queue: [S1@0]
$S1_{max}$: 0
$S2_{max}$: 0

S2

Timestamp: 0,0,0
Queue: [S1@0,0,0]

S1

Timestamp: 0,0,0
Queue: [S1@0,0,0]

S3

Timestamp: 0,0,0
Queue: [S1@0,0,0]

S2

Timestamp: 0,1,0
Queue: [S1@0,0,0]

request@0,1,0

request@0,1,0

S1

S3
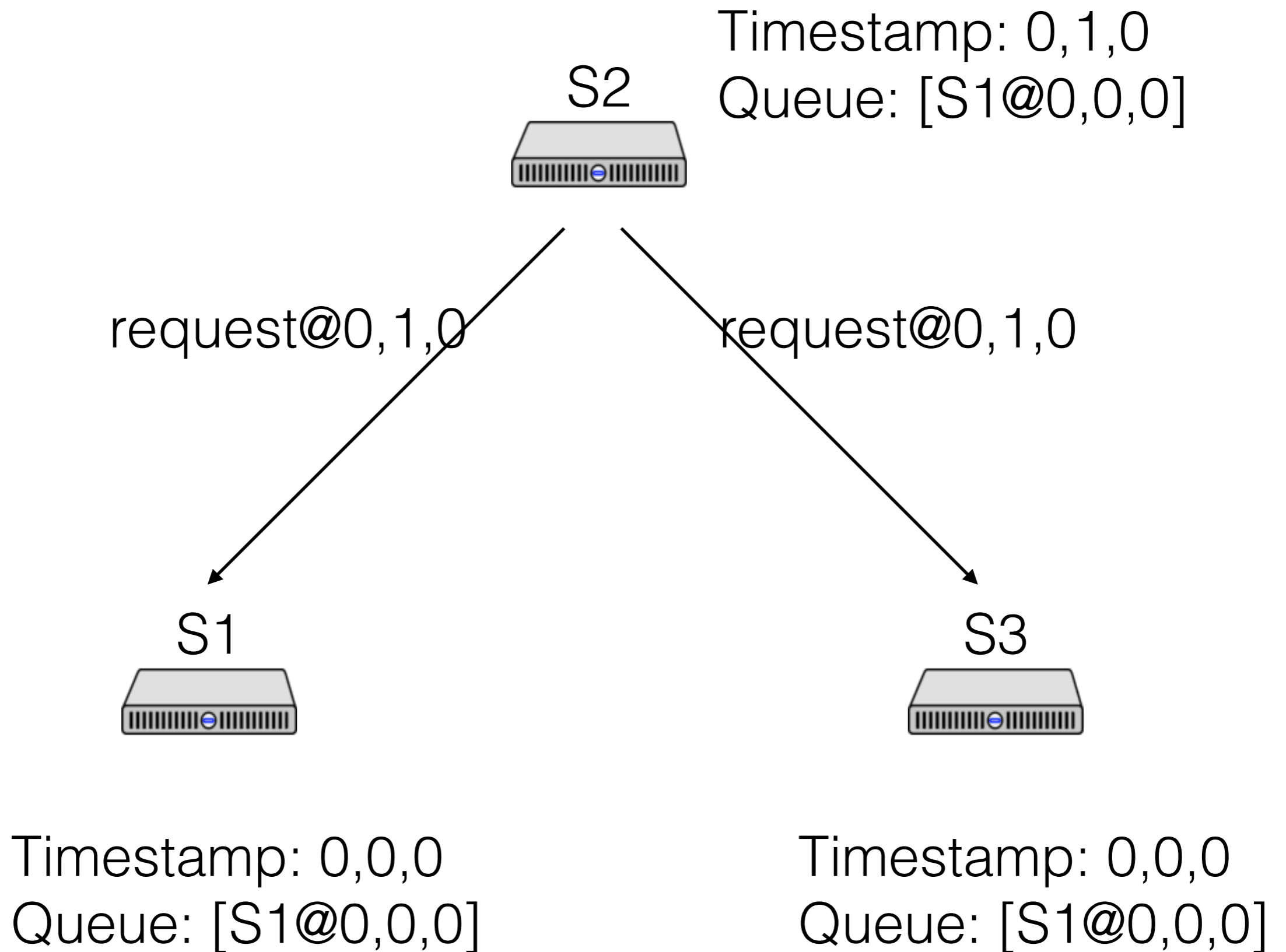
Timestamp: 0,0,0
Queue: [S1@0,0,0]

Timestamp: 0,0,0
Queue: [S1@0,0,0]

S2

Timestamp: 0,1,0
Queue: [S1@0,0,0
S2@0,1,0]

S1

Timestamp: 1,1,0
Queue: [S1@0,0,0;
S2@0,1,0]

S3

Timestamp: 0,1,1
Queue: [S1@0,0,0;
S2@0,1,0]

Timestamp: 0,1,0
Queue: [S1@0,0,0
S2@0,1,0]

S2

ack@2,1,0

ack@0,1,2

S1

S3

Timestamp: 2,1,0
Queue: [S1@0,0,0;
S2@0,1,0]

Timestamp: 0,1,2
Queue: [S1@0,0,0;
S2@0,1,0]

S2

Timestamp: 2,2,2
Queue: [S1@0,0,0
S2@0,1,0]

S1

Timestamp: 2,1,0
Queue: [S1@0,0,0;
S2@0,1,0]

S3

Timestamp: 0,1,2
Queue: [S1@0,0,0;
S2@0,1,0]

Timestamp: 2,2,2
Queue: [S1@0,0,0
S2@0,1,0]

S2

release@3,1,0

S1

release@3,1,0

S3

Timestamp: 3,1,0
Queue: [S1@0,0,0;
S2@0,1,0]

Timestamp: 0,1,2
Queue: [S1@0,0,0;
S2@0,1,0]

S2

Timestamp: 3,3,2
Queue: [S2@0,1,0]

S1

Timestamp: 3,1,0
Queue: [S2@0,1,0]

S3

Timestamp: 3,1,3
Queue: [S2@0,1,0]

S2

Timestamp: 3,4,2
Queue: [S2@0,1,0]

ack@3,4,2

S1

ack@3,1,4

S3
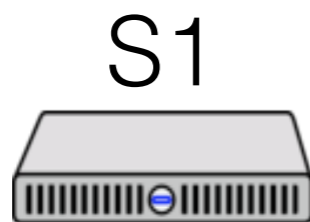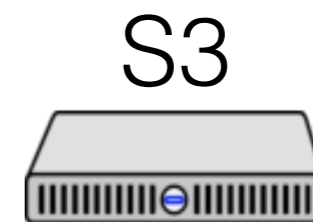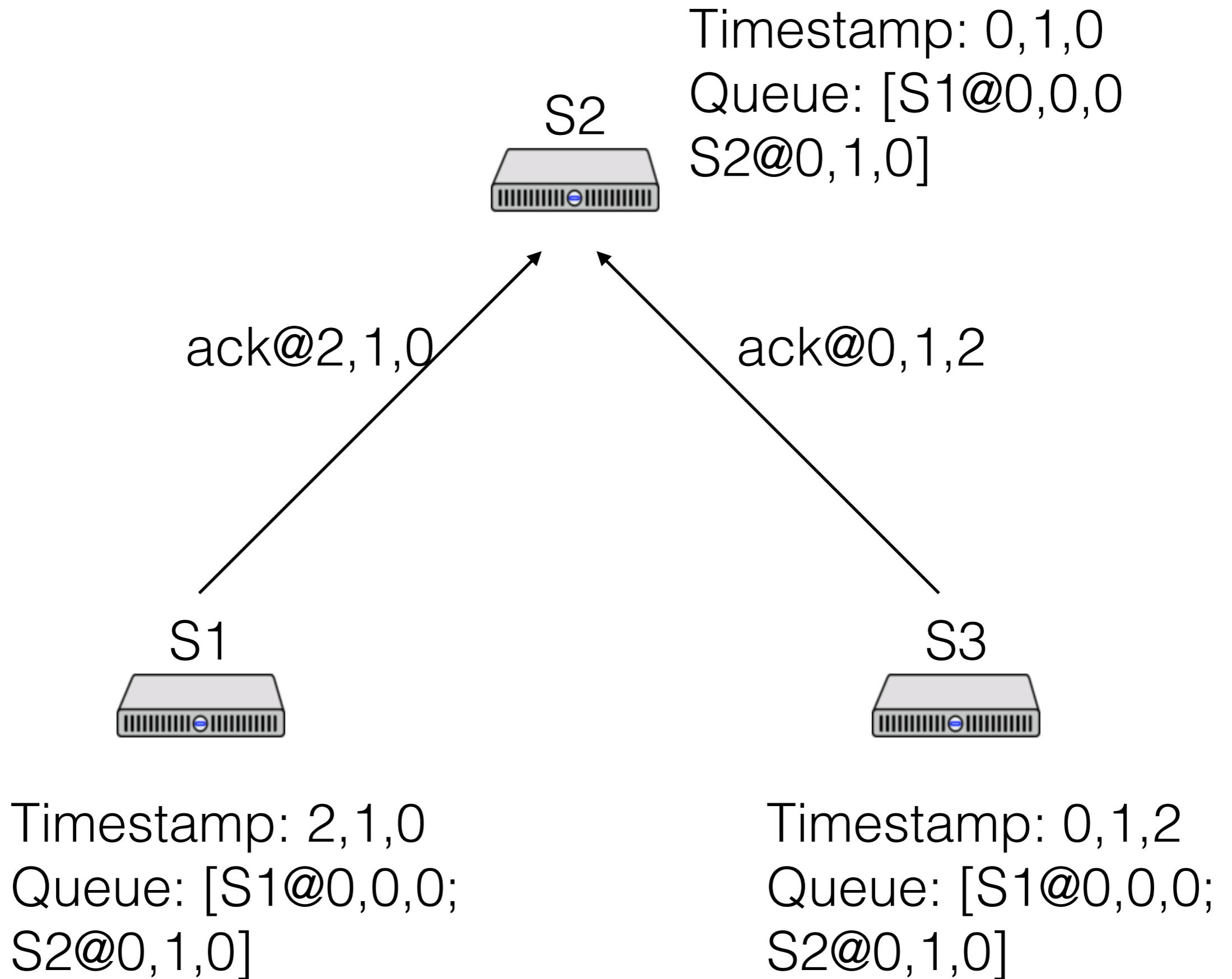
Timestamp: 3,1,0
Queue: [S2@0,1,0]

Timestamp: 3,1,4
Queue: [S2@0,1,0]

S2

Timestamp: 3,4,2
Queue: [S2@0,1,0]

S1

Timestamp: 4,4,4
Queue: [S2@0,1,0]

S3

Timestamp: 3,1,4
Queue: [S2@0,1,0]

# Some terms

Often useful: states, executions, reachability

- A state is a global state $S$ of the system: states at all nodes + channels

- An execution is a series of states $S_i$ s.t. the system is allowed to transition from $S_i$ to $S_{i+1}$

- A state $S_j$ is reachable from $S_i$ if, starting in $S_i$, it's possible for the system to end up at $S_j$

Types of properties: stable properties, invariants

- A property $P$ is stable if

$P(S_i)$ -> $P(S_{i+1})$

- A property P is an invariant if it holds on all reachable states

# Token conservation system

token

Node 1 $\longrightarrow$ Node 2

$\longleftarrow$

haveToken: bool                    haveToken: bool

In $S_o$

  - No messages

  - Node 1 has haveToken = true

  - Node 2 has haveToken = false

Nodes can send each other the token or discard the token

# Token conservation system

token

Node 1 $\longrightarrow$ Node 2

$\longleftarrow$

haveToken: bool                    haveToken: bool

Invariant: token in at most one place

Stable property: no token

# Token conservation system

token

Node 1 &rarr; Node 2

haveToken: bool                    haveToken: bool

How can we check the invariant at runtime?

How can we check the stable property at runtime?

# Distributed snapshots

Why do we want snapshots?

- Detect stable properties (e.g., deadlock)

- Distributed garbage collection

- Diagnostics (is invariant still true?)

# Distributed snapshots

Record global state of the system

- Global state: state of every node, every channel

Challenges:

- Physical clocks have skew

- State can't be an instantaneous global snapshot

- State must be consistent

# Consistent snapshots

- Consistent global state: causal dependencies are captured
  - If a snapshot of a node includes some events
    - All causally earlier events should be part of snapshots of other nodes

# Space Time Diagrams

# Cuts



A cut C is a subset of the global history of H

# Consistent Cuts

- A cut is consistent if
  - e2 is in the cut and if e1 happens before e2
    - then e1 should also be in the cut
- A consistent global state is one corresponding to a consistent cut

# Inconsistent Cut (or global state)

# Physical time algorithm

What if we could trust clocks?

Idea:

- Node: "hey, let's take a snapshot @ noon"

- At noon, everyone records state

- How to handle channels?

# Physical time algorithm

Channels:

- Timestamp all messages

- Receiver records channel state

- Channel state = messages received after noon but sent before noon

Example: is there <= 1 token in the system?

# Physical time algorithm

**11:59**

Node 1 ⟶ Node 2

haveToken = true          haveToken = false

# Physical time algorithm

**11:59**

token@11:59

Node 1 → Node 2

haveToken = false                haveToken = false

# Physical time algorithm

**12:00**

Snapshot:
- token@11:59

Node 1             Node 2

haveToken = false           haveToken = false

Snapshot:
- haveToken = false

Snapshot:
- haveToken = false

# Physical time algorithm

This seems like it works, right?

What could go wrong?

# Physical time algorithm

**11:59**                              **11:58**

Node 1 ————————————————→ Node 2

<————————————————

haveToken = true          haveToken = false

# Physical time algorithm

**12:00**                                                      **11:59**

Node 1 ——————————————————————————————→ Node 2

Node 1 ←—————————————————————————————— Node 2

haveToken = true                          haveToken = false

Snapshot:
  - haveToken = true

# Physical time algorithm

**12:00**                                    **11:59**

token@12:00

Node 1 ————————————————→ Node 2

←————————————————

haveToken = false                    haveToken = false

Snapshot:
 - haveToken = true

# Physical time algorithm

**12:00**                                                    **11:59**

Node 1 ————————————————————————→ Node 2

←————————————————————————

haveToken = false                                haveToken = true

Snapshot:
  - haveToken = true

# Physical time algorithm

**12:01**                               **12:00**

Node 1   ⟶   Node 2

haveToken = false            haveToken = true

Snapshot:
 - haveToken = true

Snapshot:
 - haveToken = true

# Avoiding inconsistencies

As we've seen, physical clocks aren't accurate enough

Need to use messages to coordinate snapshot

=> make sure Node 2 takes snapshot before receiving any messages sent after Node 1 takes snapshot

# Better algorithm

**11:59**                                    **11:58**

Node 1 ────────────────────────▶ Node 2

◀──────────────────────────

haveToken = true                haveToken = false

# Better algorithm

**12:00**                                                    **11:59**

snapshot@12:00

Node 1  →  Node 2

haveToken = true                          haveToken = false

Snapshot:
  - haveToken = true

# Better algorithm

**12:00**                                    **11:59**

token@12:00

snapshot@12:00

Node 1 ───────────────────────────────▶ Node 2

◀───────────────────────────────

haveToken = false                      haveToken = false

Snapshot:
  - haveToken = true

# Better algorithm

**12:00**                                    **11:59**

token@12:00

Node 1 ──────────────────────────▶ Node 2

          ◀──────────────────────────

haveToken = false                    haveToken = false

Snapshot:                            Snapshot:
 - haveToken = true                   - haveToken = false

# Better algorithm

**12:00**                    **11:59**

Node 1 ➝ Node 2

haveToken = false          haveToken = true

Snapshot:
  - haveToken = true

Snapshot:
  - haveToken = false

# Better algorithm

Node 1 ⟶ Node 2

Node 2 ⟶ Node 1

haveToken = false

haveToken = true

Snapshot:
  - haveToken = true

Snapshot:
  - haveToken = false

# Distributed Snapshots

As we've seen, physical clocks aren't accurate enough

Need to use messages to coordinate snapshot

=> make sure Node 2 takes snapshot before receiving any messages sent after Node 1 takes snapshot

# Chandy-Lamport Snapshots

At any time, a node can decide to snapshot

- Actually, multiple nodes can

That node:

- Records its current state

- Sends a "marker" message on all channels

When a node receives a marker, snapshot

- Record current state

- Send marker message on all channels

How to record channel state?

# Chandy-Lamport Snapshots

Channel state recorded by the receiver

Recorded when marker received on that channel

- Why do we know we'll receive a marker on every channel?

When marker received on channel, record:

- Empty, if this is the first marker

- Messages received on channel since we snapshotted, otherwise

# Chandy-Lamport Snapshots

# Chandy-Lamport Snapshots

Node 1 ⟶ Node 2

haveToken = true          haveToken = false

# Chandy-Lamport Snapshots

token

Node 1 ⟶ Node 2

haveToken = false          haveToken = false

# Chandy-Lamport Snapshots

token

Node 1 ⟶ Node 2

⟵

marker

haveToken = false          haveToken = false

Snapshot:
  - haveToken = false

# Chandy-Lamport Snapshots

marker

token

Node 1 $\longrightarrow$ Node 2

$\longleftarrow$

haveToken = false                    haveToken = false

Snapshot:                            Snapshot:
 - haveToken = false                  - haveToken = false

# Chandy-Lamport Snapshots

marker

Node 1 ──────────────────────────→ Node 2

←──────────────────────────

haveToken = false                    haveToken = true

Snapshot:                            Snapshot:
 - haveToken = false                  - haveToken = false

                                     In-flight:
                                      - token

# Chandy-Lamport Snapshots

Snapshot:
- token

Node 1 ──────────────────────────────→ Node 2

←──────────────────────────────

haveToken = false

haveToken = true

Snapshot:
 - haveToken = false

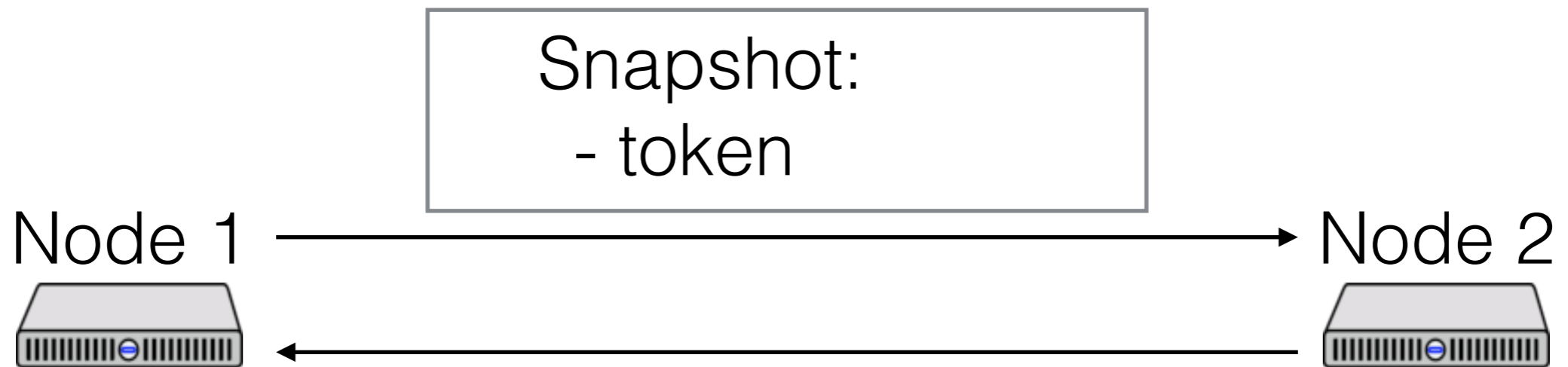Snapshot:
 - haveToken = false

# Chandy-Lamport Snapshots

What if multiple nodes initiate the snapshot?

- Follow same rules: send markers on all channels

Intuition:

- All initiators are concurrent

- Concurrent snapshots are ok, as long as we account for messages in flight

- If receive marker before initiating, must snapshot to be consistent with other nodes

# Consistent Cut

A cut is the set of events on each node in the system that are included in the snapshot

A consistent cut is a cut that respects causality

If an event is included by any node, all events that "happen before" the event are also included