# Distributed Information Processing

## 3rd Lecture

Eom, Hyeonsang (엄현상)

Department of Computer Science & Engineering

Seoul National University

# Outline

- Clock and Global States
    - Global States
    - Determining Consistent Global States
- Q&A

# Global States

- ## Prefix of Pi's History & Global History

$$h_i^k = <e_i^j \mid j = 1, ..., k>, \quad H = \bigcup_{i=1}^{N} h_i$$

- ## Cut & Frontier

$$C = \bigcup_{i=1}^{N} h_i^{C_i}, \quad F = \{e_i^{C_i} \mid i = 1, ..., N\}$$

Set of All
Affected Values

- ## Global State (Corresponding to C)

$$S = \{s_i^{C_i} \mid i = 1, ..., N, \ s_i^{C_i} \ is \ P_i's \ state \ immediately \ after \ e_i^{C_i}\}$$
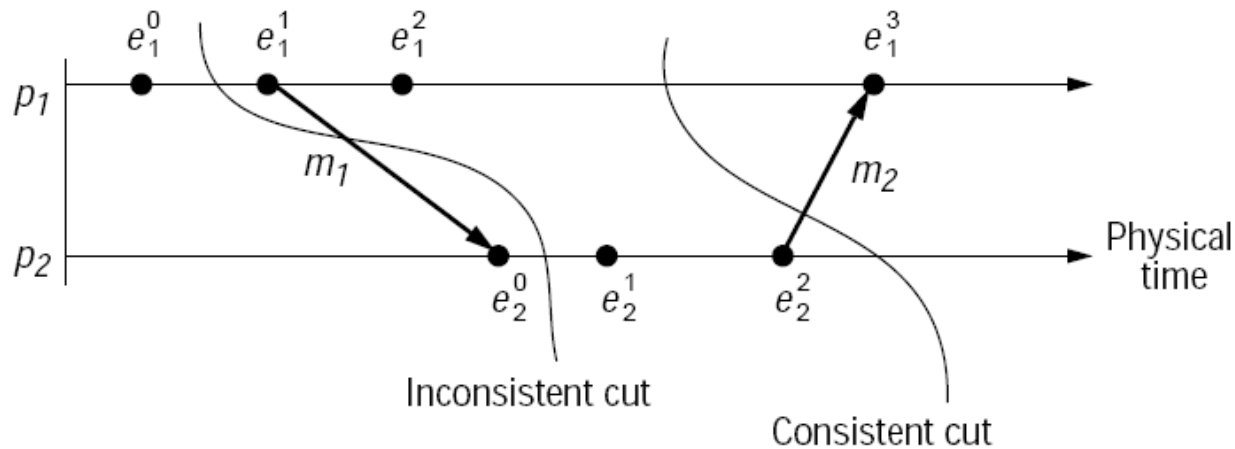
- ## Run: a Total Ordering in a Global History Consistent with Each Local History

# Consistent Cuts & Runs

■ C Is Consistent If the Following Holds:

$$\forall e \in C, \ e' \rightarrow e \ \Rightarrow \ e' \in C$$

$p_1$   $e_1^0$   $e_1^1$   $e_1^2$   $e_1^3$

$m_1$   $m_2$

Physical time

$p_2$   $e_2^0$   $e_2^1$   $e_2^2$
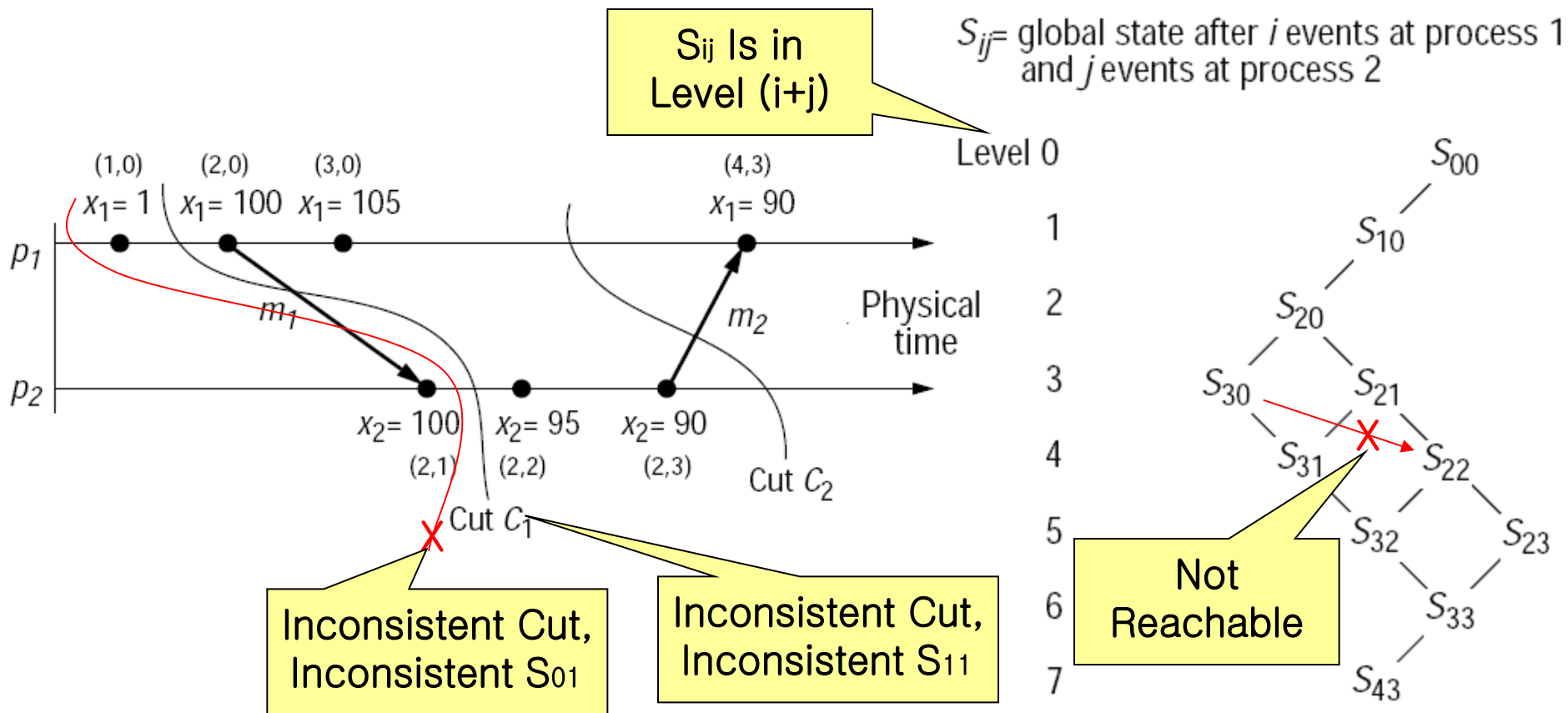
Inconsistent cut

Consistent cut

Coulouris, Dollimore and Kindberg   Distributed Systems: Concepts and Design   Edn. 4   © Pearson Education 2005

■ Consistent Run: a Total Ordering in a Consistent Global History, Consistent with the Happened-Before Relation
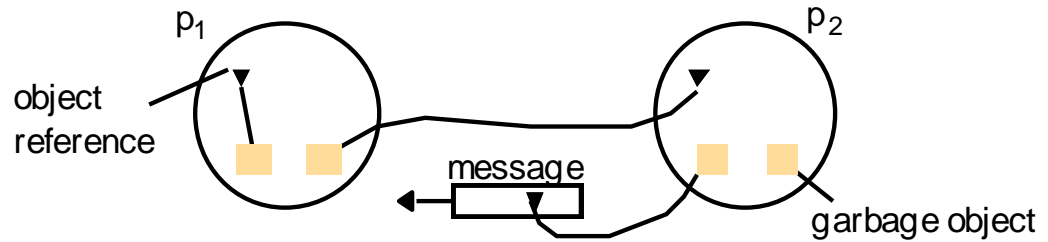
# Lattice of Global States

- Observing Consistent Global States

$$S = \{s_i \mid i = 1, ..., N\} \; Is \; Consistent \; iff \; VC_i(s_i)[i] \geq VC(s_j)[i] \; for \; i, j = 1, ..., N$$

S$_{ij}$ Is in Level (i+j)

$S_{ij}$= global state after $i$ events at process 1 and $j$ events at process 2

(1,0)    (2,0)    (3,0)                          (4,3)
$x_1$= 1   $x_1$= 100  $x_1$= 105                    $x_1$= 90

$p_1$

$m_1$                                             $m_2$                Physical time

$p_2$

$x_2$= 100   $x_2$= 95   $x_2$= 90                      Cut $C_2$
(2,1)      (2,2)      (2,3)

Cut $C_1$

Inconsistent Cut, Inconsistent S$_{01}$

Inconsistent Cut, Inconsistent S$_{11}$

Not Reachable

Level 0    $S_{00}$
1          $S_{10}$
2          $S_{20}$
3          $S_{30}$    $S_{21}$
4          $S_{31}$    $S_{22}$
5          $S_{32}$    $S_{23}$
6          $S_{33}$
7          $S_{43}$

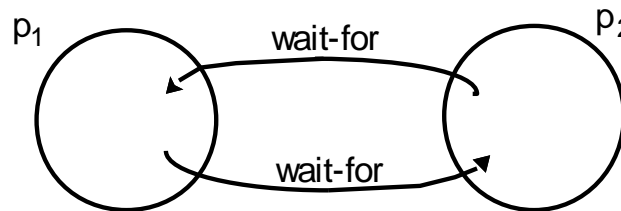Coulouris, Dollimore and Kindberg   Distributed Systems: Concepts and Design   Edn. 4   © Pearson Education 2005

# Detecting Global Properties



a. Garbage collection

b. Deadlock

c. Termination

# Distributed 'Snapshot' Algorithm [Chandy85]

- ## Consistent Global-State Detection

*Marker sending rule for process pi*

**Marker to Record the State**

> After $p_i$ has recorded its state, for each outgoing channel $c$:
>
> > $p_i$ sends one marker message over $c$
> >
> > (before it sends any other message over $c$).

*Marker receiving rule for process $p_i$*

> On $p_i$'s receipt of a *marker* message over channel $c$:
>
> *if* ($p_i$ has not yet recorded its state) it
>
> > records its process state now;
> >
> > records the state of $c$ as the empty set;
> >
> > turns on recording of messages arriving over other incoming channels;
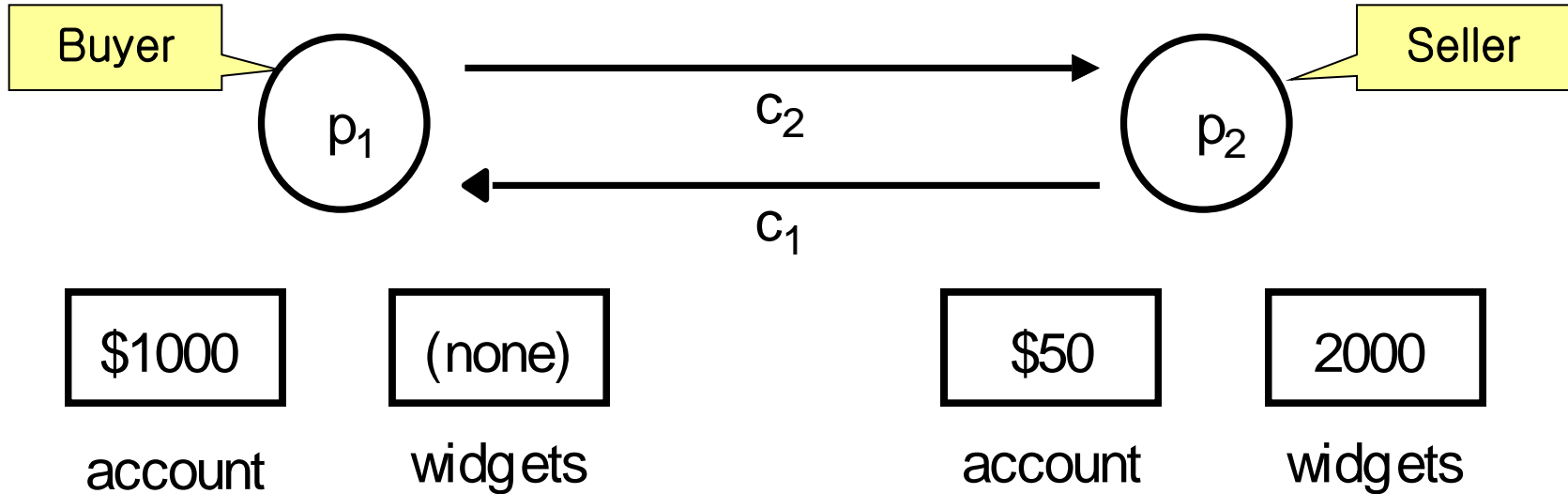>
> *else*
>
> > $p_i$ records the state of $c$ as the set of messages it has received over $c$
> > since it saved its state.
>
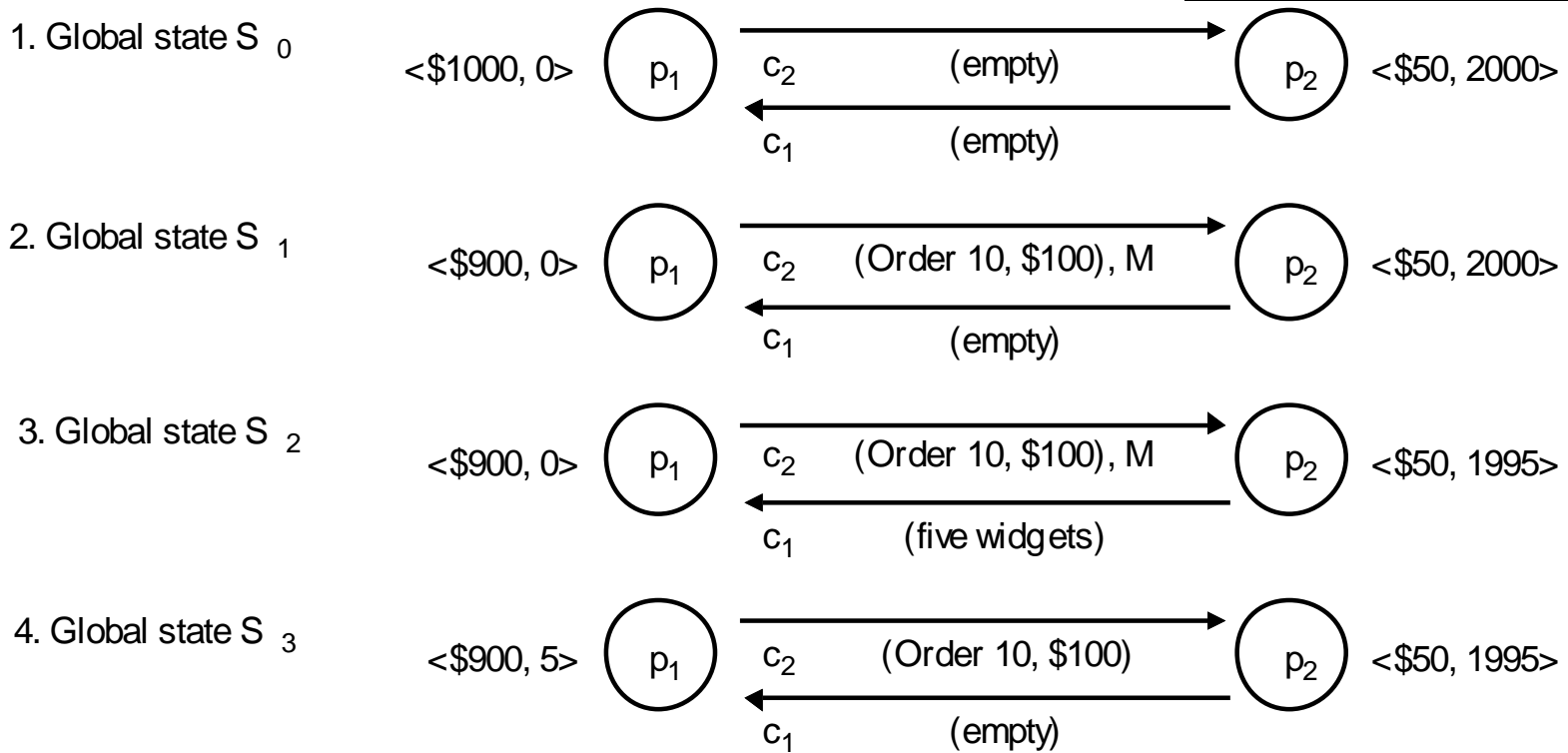> *end if*

# Illustration: How the Alg. Works

- Initial States of the Components

# Illustration (Cont'd)

1. Global state $S_0$

$\langle\$1000, 0\rangle$   $p_1$   $c_2$   (empty)   $p_2$   $\langle\$50, 2000\rangle$

$c_1$   (empty)

2. Global state $S_1$

$\langle\$900, 0\rangle$   $p_1$   $c_2$   (Order 10, \$100), M   $p_2$   $\langle\$50, 2000\rangle$

$c_1$   (empty)

3. Global state $S_2$

$\langle\$900, 0\rangle$   $p_1$   $c_2$   (Order 10, \$100), M   $p_2$   $\langle\$50, 1995\rangle$

$c_1$   (five widgets)

4. Global state $S_3$

$\langle\$900, 5\rangle$   $p_1$   $c_2$   (Order 10, \$100)   $p_2$   $\langle\$50, 1995\rangle$

$c_1$   (empty)
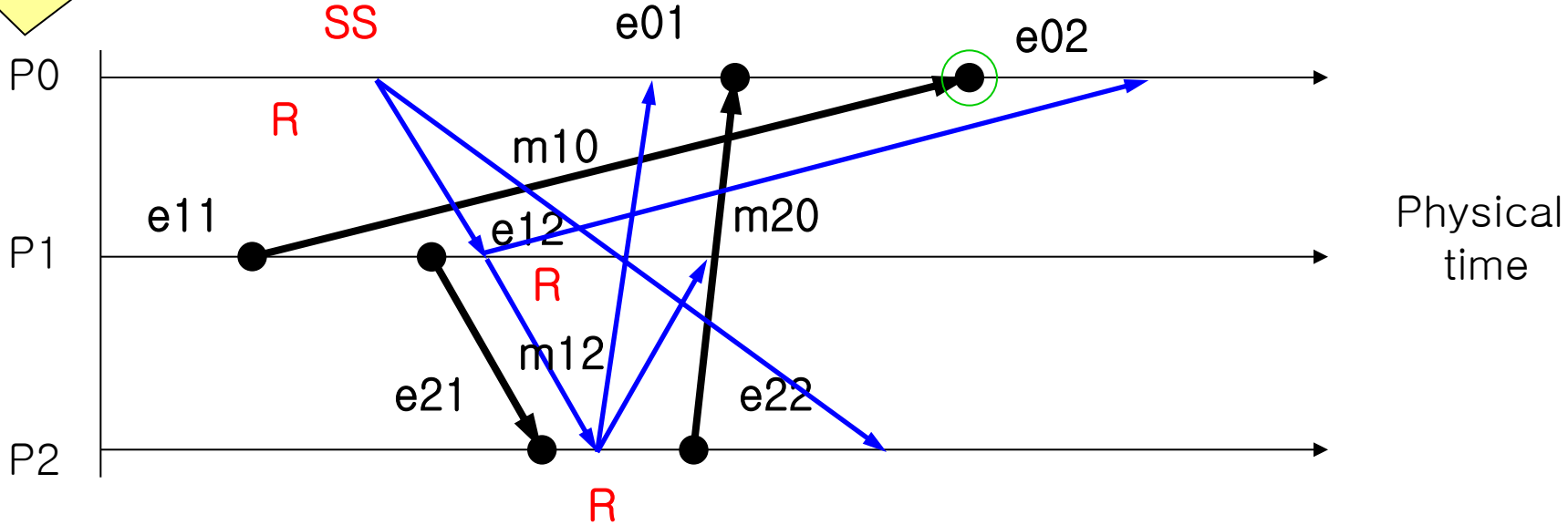
(M = marker message)

- $S = p1: \langle\$1000,0\rangle; p2:\langle\$50,1995\rangle;$
  $c1: \langle(\text{five widgets})\rangle; c2:\langle\rangle$

# Illustration w/ a Diagram



■ SS: $p0$: <>; $p1$: <$e11,e12$>; $p2$:<$e21$>

     $c01$:<>; $c02$ <>; $c10$<$m10$>; $c12$ <>

     $c20$ <>; $c21$<>

# Consistency Proof

- States Recorded by the Alg. Are Consistent:

$$\forall e_j \in C, \ e_i \rightarrow e_j \ \Rightarrow \ e_i \in C$$

$Show: \ e_i \notin C, \ e_i \rightarrow e_j \ \Rightarrow \ e_j \notin C$ — i≠j

- Assume That $P_i$ Recorded Its State before $e_i$
- Marker Would Have Reached $P_j$ before the Message for $e_j$
- $P_j$ Would Have Recorded Its State before $e_j$

# Characteristics of Snapshots

- Derivation of "Observed" Run from "Actual" Run

$$"Actual"\ Run : < e_i^k \mid i = 1, ..., N >=< e^j >$$

$$"Permuted"\ Run : < ..., e^{R-1}, e^R, ... >$$

Recodring

Observed
(Consistent ) Run

  - A Non-Observed Event May Occur before an Observed Event in the "Actual" Run

  - If a Non-Observed Event Precedes an Observed Event (Next to it) in the "Actual" Run, the Events Can Be Swapped Preserving Consistency

# Global State Predicates

- **Functions That Map Global States to True or False**
  - Stable: Once True, It Remains True
    - E.g., deadlock or termination
  - Unstable: Not Stable
    - Possibly True: True At Some Point
      - E.g., snapshot by the 'Snapshot' Algorithm
    - Definitely True: True in All Cases