

## 2.1 Asynchronous execution with synchronous communication

When all the communication between pairs of processes is by using synchronous send and receive primitives, the resulting order is synchronous order. The algorithms run on asynchronous systems will not work in synchronous system and vice versa is also true.

### Realizable Synchronous Communication (RSC)

*A-execution can be realized under synchronous communication is called a realizable with synchronous communication (RSC).*

- An execution can be modeled to give a total order that extends the partial order  $(E, <)$ .
- In an A-execution, the messages can be made to appear instantaneous if there exist a linear extension of the execution, such that each send event is immediately followed by its corresponding receive event in this linear extension.

*Non-separated linear extension is an extension of  $(E, <)$  is a linear extension of  $(E, <)$  such that for each pair  $(s, r) \in T$ , the interval  $\{x \in E \mid s < x < r\}$  is empty.*

*A A-execution  $(E, <)$  is an RSC execution if and only if there exists a non-separated linear extension of the partial order  $(E, <)$ .*

- In the non-separated linear extension, if the adjacent send event and its corresponding receive event are viewed atomically, then that pair of events shares a common past and a common future with each other.

### Crown

*Let  $E$  be an execution. A crown of size  $k$  in  $E$  is a sequence  $\langle (s^i, r^i), i \in \{0, \dots, k-1\} \rangle$  of pairs of corresponding send and receive events such that:  $s^0 < r^1, s^1 < r^2, s^{k-2} < r^{k-1}, s^{k-1} < r^0$ .*

The crown is  $\langle (s^1, r^1) (s^2, r^2) \rangle$  as we have  $s^1 < r^2$  and  $s^2 < r^1$ . Cyclic dependencies may exist in a crown. The crown criterion states that an A-computation is RSC, i.e., it can be realized on a system with synchronous communication, if and only if it contains no crown.

### Timestamp criterion for RSC execution

An execution  $(E, <)$  is RSC if and only if there exists a mapping from  $E$  to  $T$  (scalar timestamps) such that

- for any message  $M$ ,  $T(s(M)) = T(r(M))$ ;
- for each  $(a, b)$  in  $(E \times E) \setminus T$ ,  $a < b \implies T(a) < T(b)$

#### 2.2.1 Hierarchy of ordering paradigms

The orders of executions are:

- Synchronous order (SYNC)
- Causal order (CO)
- FIFO order (FIFO)
- Non FIFO order (non-FIFO)

#### The Execution order have the following results

- For an A-execution, A is RSC if and only if A is an S-execution.
- $RSC \subset CO \subset FIFO \subset A$
- This hierarchy is illustrated in Figure 2.3(a), and example executions of each class are shown side-by-side in Figure 2.3(b)

- The above hierarchy implies that some executions belonging to a class X will not belong to any of the classes included in X. The degree of concurrency is most in A and least in SYNC.
- A program using synchronous communication is easiest to develop and verify.
- A program using non-FIFO communication, resulting in an A execution, is hardest to design and verify.

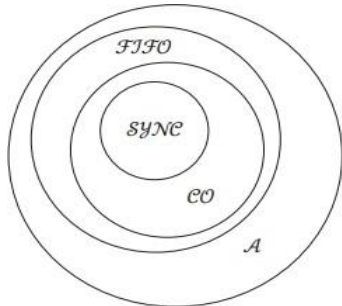


Fig (a)

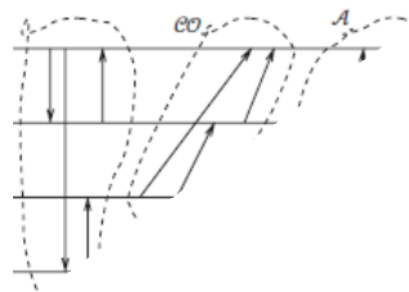


Fig (b)

Fig 2.3: Hierarchy of execution classes

### 2.2.3 Simulations

- The events in the RSC execution are scheduled as per some non-separated linear extension, and adjacent (s, r) events in this linear extension are executed sequentially in the synchronous system.
- The partial order of the asynchronous execution remains unchanged.
- If an A-execution is not RSC, then there is no way to schedule the events to make them RSC, without actually altering the partial order of the given A-execution.
- However, the following indirect strategy that does not alter the partial order can be used.
- Each channel  $C_{i,j}$  is modeled by a control process  $P_{i,j}$  that simulates the channel buffer.
- An asynchronous communication from  $i$  to  $j$  becomes a synchronous communication from  $i$  to  $P_{i,j}$  followed by a synchronous communication from  $P_{i,j}$  to  $j$ .
- This enables the decoupling of the sender from the receiver, a feature that is essential in asynchronous systems.

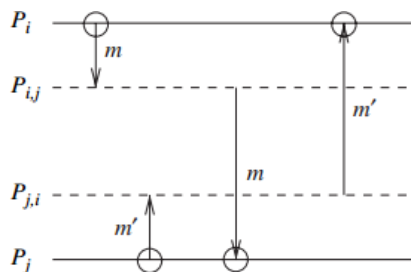


Fig 2.4: Modeling channels as processes to simulate an execution using asynchronous primitives on synchronous system

### Synchronous programs on asynchronous systems

- A (valid) S-execution can be trivially realized on an asynchronous system by scheduling the messages in the order in which they appear in the S-execution.
- The partial order of the S-execution remains unchanged but the communication occurs

on an asynchronous system that uses asynchronous communication primitives.

- Once a message send event is scheduled, the middleware layer waits for acknowledgment; after the ack is received, the synchronous send primitive completes.

## 2.2 SYNCHRONOUS PROGRAM ORDER ON AN ASYNCHRONOUS SYSTEM

### Non deterministic programs

The partial ordering of messages in the distributed systems makes the repeated runs of the same program will produce the same partial order, thus preserving deterministic nature. But sometimes the distributed systems exhibit non determinism:

- A receive call can receive a message from any sender who has sent a message, if the expected sender is not specified.
- Multiple send and receive calls which are enabled at a process can be executed in an interchangeable order.
- If  $i$  sends to  $j$ , and  $j$  sends to  $i$  concurrently using blocking synchronous calls, there results a deadlock.
- There is no semantic dependency between the send and the immediately following receive at each of the processes. If the receive call at one of the processes can be scheduled before the send call, then there is no deadlock.

### 2.3.1 Rendezvous

Rendezvous systems are a form of synchronous communication among an arbitrary number of asynchronous processes. All the processes involved meet with each other, i.e., communicate synchronously with each other at one time. Two types of rendezvous systems are possible:

- Binary rendezvous: When two processes agree to synchronize.
- Multi-way rendezvous: When more than two processes agree to synchronize.

### Features of binary rendezvous:

- For the receive command, the sender must be specified. However, multiple receive commands can exist. A type check on the data is implicitly performed.
- Send and received commands may be individually disabled or enabled. A command is disabled if it is guarded and the guard evaluates to false. The guard would likely contain an expression on some local variables.
- Synchronous communication is implemented by scheduling messages under the covers using asynchronous communication.
- Scheduling involves pairing of matching send and receives commands that are both enabled. The communication events for the control messages under the covers do not alter the partial order of the execution.

### 2.3.2 Binary rendezvous algorithm

If multiple interactions are enabled, a process chooses one of them and tries to synchronize with the partner process. The problem reduces to one of scheduling messages satisfying the following constraints:

- Schedule on-line, atomically, and in a distributed manner.
- Schedule in a deadlock-free manner (i.e., crown-free).
- Schedule to satisfy the progress property in addition to the safety property.

**Steps in Bagrodia algorithm**

1. Receive commands are forever enabled from all processes.
2. A send command, once enabled, remains enabled until it completes, i.e., it is not possible that a send command gets before the send is executed.
3. To prevent deadlock, process identifiers are used to introduce asymmetry to break potential crowns that arise.
4. Each process attempts to schedule only one send event at any time.

The message (M) types used are: M, ack(M), request(M), and permission(M). Execution events in the synchronous execution are only the send of the message M and receive of the message M. The send and receive events for the other message types – ack(M), request(M), and permission(M) which are control messages. The messages request(M), ack(M), and permission(M) use M's unique tag; the message M is not included in these messages.

---

(message types)

M, ack(M), request(M), permission(M)

**(1)  $P_i$  wants to execute SEND(M) to a lower priority process  $P_j$ :**

$P_i$  executes *send*(M) and blocks until it receives *ack*(M) from  $P_j$ . The send event SEND(M) now completes.

Any  $M'$  message (from a higher priority processes) and *request*( $M'$ ) request for synchronization (from a lower priority processes) received during the blocking period are queued.

**(2)  $P_i$  wants to execute SEND(M) to a higher priority process  $P_j$ :**

(2a)  $P_i$  seeks permission from  $P_j$  by executing *send*(request(M)).

// to avoid deadlock in which cyclically blocked processes queue // messages.

(2b) While  $P_i$  is waiting for permission, it remains unblocked.

(i) If a message  $M'$  arrives from a higher priority process  $P_k$ ,  $P_i$  accepts  $M'$  by scheduling a RECEIVE( $M'$ ) event and then executes *send*(ack( $M'$ )) to  $P_k$ .

(ii) If a *request*( $M'$ ) arrives from a lower priority process  $P_k$ ,  $P_i$  executes *send*(permission( $M'$ )) to  $P_k$  and blocks waiting for the message  $M'$ . When  $M'$  arrives, the RECEIVE( $M'$ ) event is executed.

(2c) When the *permission*(M) arrives,  $P_i$  knows partner  $P_j$  is synchronized and  $P_i$  executes *send*(M). The SEND(M) now completes.

**(3) request(M) arrival at  $P_i$  from a lower priority process  $P_j$ :**

At the time a *request*(M) is processed by  $P_i$ , process  $P_i$  executes *send*(permission(M)) to  $P_j$  and blocks waiting for the message M. When M arrives, the RECEIVE(M) event is executed and the process unblocks.

**(4) Message M arrival at  $P_i$  from a higher priority process  $P_j$ :**

At the time a message M is processed by  $P_i$ , process  $P_i$  executes RECEIVE(M) (which is assumed to be always enabled) and then *send*(ack(M)) to  $P_j$ .

**(5) Processing when  $P_i$  is unblocked:**

When  $P_i$  is unblocked, it dequeues the next (if any) message from the queue and processes it as a message arrival (as per rules 3 or 4).

-----

**Fig 2.5: Bagrodia Algorithm**

## 2.3 GROUP COMMUNICATION

Group communication is done by broadcasting of messages. A message broadcast is the sending of a message to all members in the distributed system. The communication may be

- **Multicast:** A message is sent to a certain subset or a group.
- **Unicasting:** A point-to-point message communication.

The network layer protocol cannot provide the following functionalities:

- Application-specific ordering semantics on the order of delivery of messages.
- Adapting groups to dynamically changing membership.
- Sending multicasts to an arbitrary set of processes at each send event.
- Providing various fault-tolerance semantics.
- The multicast algorithms can be open or closed group.

### Differences between closed and open group algorithms:

Closed group algorithms	Open group algorithms
If sender is also one of the receiver in the multicast algorithm, then it is closed group algorithm.	If sender is not a part of the communication group, then it is open group algorithm.
They are specific and easy to implement.	They are more general, difficult to design and expensive.
It does not support large systems where client processes have short life.	It can support large systems.

## 2.4 CAUSAL ORDER (CO)

In the context of group communication, there are two modes of communication: *causal order and total order*. Given a system with FIFO channels, causal order needs to be explicitly enforced by a protocol. The following two criteria must be met by a causal ordering protocol:

- **Safety:** In order to prevent causal order from being violated, a message  $M$  that arrives at a process may need to be buffered until all system wide messages sent in the causal past of the send ( $M$ ) event to that same destination have already arrived. The arrival of a message is transparent to the application process. The delivery event corresponds to the receive event in the execution model.
- **Liveness:** A message that arrives at a process must eventually be delivered to the process.

### 2.5.1 The Raynal–Schiper–Toueg algorithm

- Each message  $M$  should carry a log of all other messages sent causally before  $M$ 's send event, and sent to the same destination  $\text{dest}(M)$ .
- The Raynal–Schiper–Toueg algorithm canonical algorithm is a representative of several algorithms that reduces the size of the local space and message space overhead by various techniques.

- This log can then be examined to ensure whether it is safe to deliver a message.
- All algorithms aim to reduce this log overhead, and the space and time overhead of maintaining the log information at the processes.
- To distribute this log information, broadcast and multicast communication is used.
- The hardware-assisted or network layer protocol assisted multicast cannot efficiently provide features:
  - Application-specific ordering semantics on the order of delivery of messages.
  - Adapting groups to dynamically changing membership.
  - Sending multicasts to an arbitrary set of processes at each send event.
  - Providing various fault-tolerance semantics