

Nicolò Rivetti

Concurrent and Parallel Programming Project

# Total Order Broadcast

# Project Description

This project aims to implement, using **C/MPI**, the **Total Order Broadcast** abstraction presented in the book “Introduction to Reliable and Secure Distributed Programming” by Cachin, Guerraoui and Rodrigues [1].

The abstraction is built through the following **protocol stack**:

- **Consensus Based Total Order Broadcast**
- **Flooding Consensus**
- **Eager Reliable Broadcast**
- **Best Effort Broadcast**
- **Perfect Failure Detector**
- **Perfect Point to Point Link**

# Building Blocks

The algorithms are based on a **asynchronous event based composition model**, hence each layer is modeled as components which interact through events generation and consumption.

- **Event Queue & Runtime**

The protocol layers **register events** to the first, while the latter **consumes events** causing the related handlers to be executed.

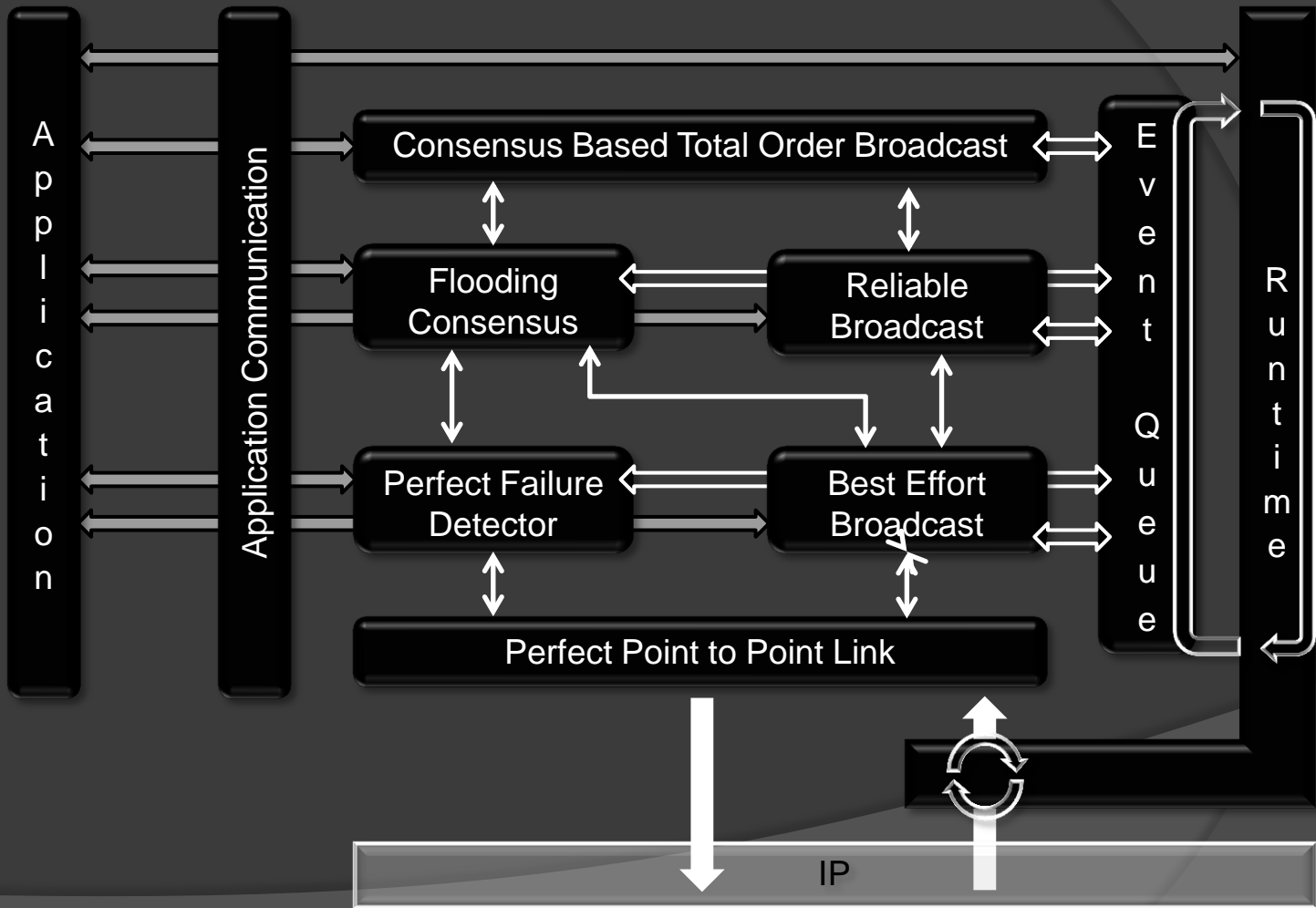
- **Perfect P2P Link & Runtime**

MPI provides a point to point communication primitive fitting the PP2P Link properties, however the messages reception is not completely asynchronous. To add this capability, the runtime takes care of checking for new messages and post them as event to the relevant component.

- **Application Communication**

To allow maximal freedom to the application layer, the library lives in it's own process, hence this module hides the communication details between the runtime and the application process.

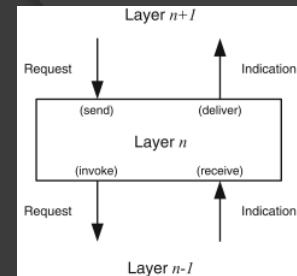
# Overall Architecture



# Event Queue

Using a event based composition model, all module share a common structure which implements the **Indication / Request** API shown in [1].

- Each module allows upper layers to **register a callback** which is a function pointer to a procedure that will be called as **Indication**.
- Each module expose a procedure to be called as **Request**.
- The Request and Indication procedures implementation will only **register an event** to the Event Queue, as a function pointer, a pointer to data and the data size



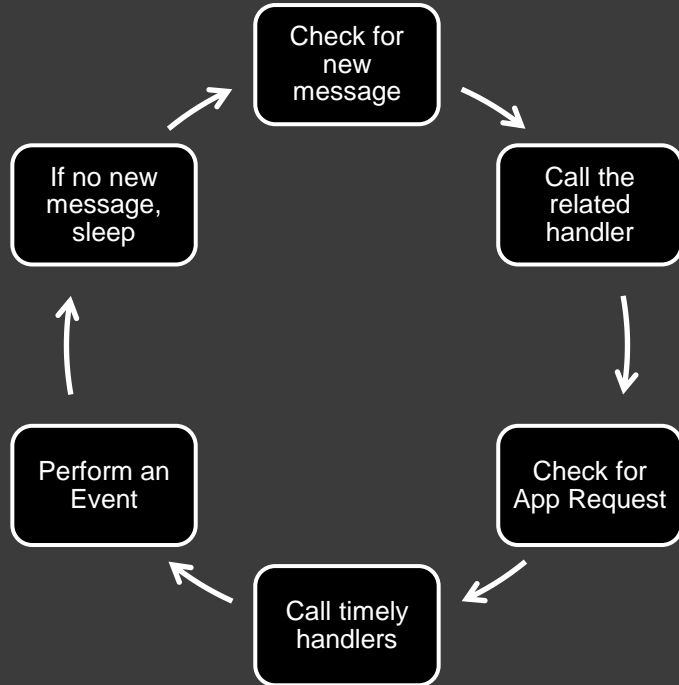
```
20 typedef struct event_struct {
21     int (*callback)(int * data, int size);
22     void * data;
23     int size;
24     struct event_struct *next;
25 } event_t;
```

```
116 int _perform_event() {
117
118     // There are no events to perform
119     if (event_head->next == NULL ) {
120         return 0;
121     }
122
123     sprintf(log_buffer, "RUNTIME INFO: Performing Event (Events in queue:%d).", event_num);
124     runtime_log(DEBUG_EVENTQUEUE);
125
126     event_t * event = event_head->next;
127
128     // Performing and removing the event
129     int ret = event->callback(event->data, event->size);
130
131     // If it's the last event in the queue, event_tail must point to the head
132     if (event_tail == event_head->next) {
133         event_tail = event_head;
134         event_head->next = NULL;
135     } else {
136         event_head->next = event->next;
137     }
138
139     // Deallocating event_struct, pointed data should be handled by callback
140     free(event);
141
142     event_num--;
143
144     return ret;
145 }
```

```
26 // Abstraction Logic Implementation
27 int _best_effort_broadcast_send(int *message, int count, int handle);
28 int _best_effort_broadcast_send_callback(int *data, int size);
29 int _best_effort_broadcast_receive(int *message, int count, int sender);
30
31 // Commodity and Utility Functions
32 int _best_effort_broadcast_register_callback(int (*callback_func)(int*, int, int), int* handle);
```

# Runtime

The Runtime is the heart of the system. After setting up all modules required by the application (using flags), it loops performing the following tasks:



- Through the `MPI_Iprobe` function, the runtime checks if a new messages is ready to be received.
- If it's the case, then it actually receives the message. Given the TAG of the message, the callback registered whit it is called
- Checks if the application layer has sent some messages through the App Comm Module
- Execute all the callback registered as timely.
- Consume one event registered on the Queue Event
- If no messages has been received, it sleeps for few ms

# Serial Number Management (1/2)

Many protocols are based on the assumption that each message is tagged with a **unique identifier**, which could be the pair <senderId, Serial Number>. Eventually, since **memory is limited**, the Serial Number (SN) will **overflow** and must be set back to 0, **breaking the uniqueness assumption**.

In this project two way of dealing with this issue have been devised, the first is used to manage the *Uniform Consensus* Instances SN

```
843  /*
844  * Commodity function to retrieve the instance struct (if it exists) given it's id
845  */
846  instance * getinstance(int id) {
847      instance *parent = head;
848
849      gettimeofday(&gettimeofdayArg, NULL );
850      unsigned long currentTS = gettimeofdayArg.tv_sec;
851
852      while (parent->next != NULL ) {
853          if ((parent->next->decision != NULL ) && (parent->next->decidedTs < currentTS - MIN_DELAY)) {
854              instance *ist = parent->next;
855
856              // Delete Record
857
858              parent->next = parent->next->next;
859              free(ist);
860          } else {
861              if (parent->next->instance_id == id) {
862                  return parent->next;
863              }
864              parent = parent->next;
865          }
866      }
867
868      return NULL ;
869  }
870  }
```

```
79  struct instance_struct {
80      int instance_id;
81      int callback;
82      int upper_layer_id;
83      int round;
84      proposal **proposals;
85      proposal *decision;
86      unsigned long decidedTs;
87      int **receivedFrom;
88      struct instance_struct *next;
89  };
90
91  typedef struct instance_struct instance;
92  static instance* head;
```

- A new SN is accepted only if it has not been recorded in the SN data structure.
- When a SN is stored, the record contains also a timestamp (TS).
- Each time the data structure is looked up, if record TS < current TS - Delay, then the record is deleted

# Serial Number Management (2/2)

This mechanism is used to manage the arriving messages for *Eager Reliable Broadcast* and *Consensus Based Total Order Broadcast*

The second solution is **less memory consuming** but a bit more complex.

The available space is divided in quarters, the four values are stored in a data structure as G0, G1, G2, G3 (Guards) and have a TS associated. T is the next expected SN.

```
11 struct serialNumber_struct {
12     unsigned int threshold;
13     unsigned int max;
14     unsigned int guard;
15     unsigned int guards[4];
16     long timestamps[4];
17 };

13 #define MAX_UINT ((unsigned int) ~0)
14 #define QUARTER0 ((unsigned int) 0)
15 #define QUARTER1 ((unsigned int) (MAX_UINT - 1) / 4)
16 #define QUARTER2 ((unsigned int) QUARTER1 * 2 + 1)
17 #define QUARTER3 ((unsigned int) QUARTER1 * 3 + 2)

78 int incrementNextSerialNumber(serialNumber *sn, unsigned int id, long ts) {
79     if (sn->threshold < sn->guards[sn->guard]) {
80         if (id >= sn->threshold && id < sn->guards[sn->guard]) {
81             if (id == sn->threshold) {
82                 incrHighThreshold(sn);
83                 incrLowThreshold(sn, ts);
84             }
85             return 0;
86         }
87     } else if (sn->threshold > sn->guards[sn->guard]) {
88         if ((id >= sn->threshold && id <= sn->max) || (id >= 0 && id < sn->guards[sn->guard])) {
89             if (id == sn->threshold) {
90                 incrHighThreshold(sn);
91                 incrLowThreshold(sn, ts);
92             }
93             return 0;
94         }
95     }
96     printf("ERROR! threshold: %u, id: %u, ts: %u, guard: %u, guardValue: %u\n", sn->threshold, id, ts, sn->guard,
97           sn->guards[sn->guard]);
98     return -1;
99 }
```

- A new SN is accepted only if  $T \leq SN < G$ , and has not been recorded in the SN data structure.
- When a SN is stored, if  $SN = T$ , then T is increased. If  $SN > \text{current } G$  and  $TS > \text{current } G \text{ TS}$ , then G is moved to the next Guard.
- Each time the data structure is looked up, all  $SN = T$  are deleted (T is increased after each deletion)

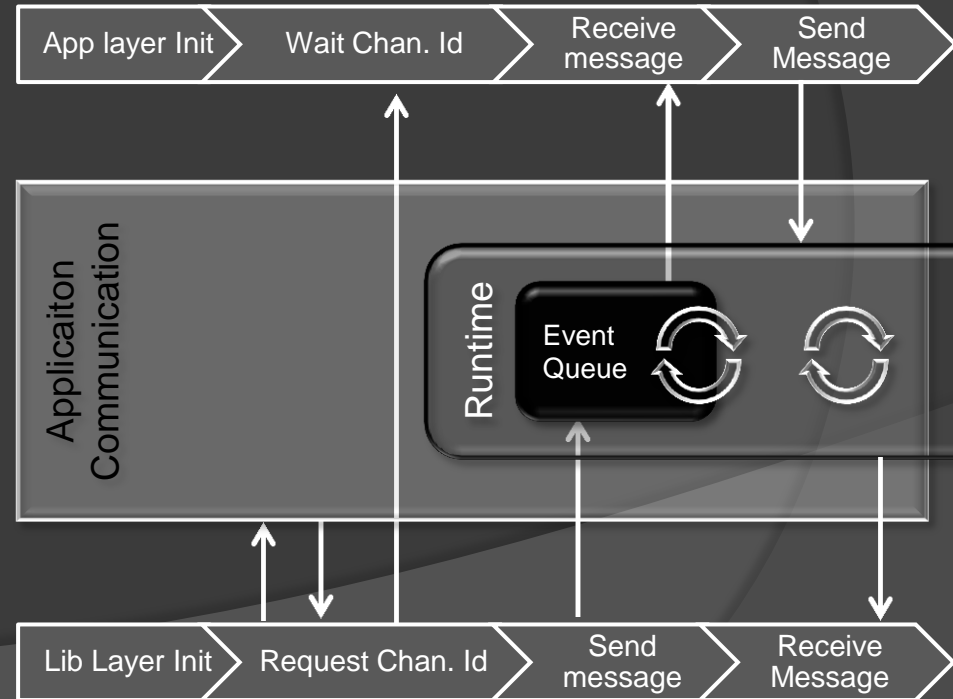


# Application Communication

Only one instance of MPI can be created per process. To **avoid unwanted interaction** between MPI API calls either the application layer must not use MPI or the runtime must live in **a new process**.

The latter was considered the best approach, requiring a module managing inter process communication

- Each module can create a channel through the App Comm module with the Application layer.
- The communication is done using **message queues**. Channels are implemented with **messages tag**.
- The initialization function of a module will ask the App Comm module for a free channel Id.
- The two parties of each module, given the channel Id can use the App Comm API to communicate



# Perfect Failure Detector

- PFD1: Strong completeness: Eventually, every process that crashes is permanently detected by every correct process.
- PFD2: Strong accuracy: If a process  $p$  is detected by any process, then  $p$  has crashed.

The Perfect failure detector requires is made of two main functions. One is registered as timely callback and the other as receive callback to the runtime:

- One check if all currently correct process have replied to the Heartbeat Request and trigger a crash event if not, calling the handler registered from the other layers and flagging the process as crashed.
- The other is triggered by received messages, replies to Heartbeat Request or registers Heartbeat Replies.

```
332 // Check for each process
333 for (i = 0; i < num_procs; i++) {
334
335     if ((detected[i] == 0) && (alive[i] == 0)) {
336         // The process has not replied to the heartbeat request and was not already detected,
337         // hence the process is flagged as failed and is added to the process crashed in this
338         // period
339         sprintf(log_buffer, "FD INFO: Failure of process %d", i);
340         runtime_log(INFO_PERFECTFAILUREDETECTOR);
341         detected[i] = 1;
342         detected_num++;
343         currently_detected[currently_detected_num] = i;
344         currently_detected_num++;
345     } else if (detected[i] == 0) {
346         // If the process was not detected and has replied to the heartbeat request,
347         // then the process is alive and a new heartbeat request is issued
348         sprintf(log_buffer, "FD INFO: Process %d alive", i);
349         runtime_log(INFO_PERFECTFAILUREDETECTOR);
350     }
351     #ifdef BUFFEREDSEND
352     MPI_Bsend(&heartbeat_request_message, 1, MPI_INT, i, runtime_receive_handle, comm);
353     #endif
354 }
```

```
upon event { Timeout } do
  forall p ∈ Π do
    if (p ∉ alive) ∧ (p ∉ detected) then
      detected := detected ∪ {p};
      trigger { P, Crash | p };
      trigger { pl, Send | p, [HEARTBEATREQUEST] };
      alive := ∅;
      starttimer(Δ);

upon event { pl, Deliver | q, [HEARTBEATREQUEST] } do
  trigger { pl, Send | q, [HEARTBEATREPLY] };

upon event { pl, Deliver | p, [HEARTBEATREPLY] } do
  alive := alive ∪ {p};
```

The buffered mode send (MPI\_Bsend) is used in this module since a fixed size message is exchanged. Hence, given the number of processes it is easy to compute an upper bound on the buffer size to be attached in order to allow the buffered communication.

# Best Effort Broadcast

- BEB1: Validity: If a correct process broadcasts a message  $m$ , then every correct process eventually delivers  $m$ .
- BEB2: No duplication: No message is delivered more than once.
- BEB3: No creation: If a process delivers a message  $m$  with sender  $s$ , then  $m$  was previously broadcast by process  $s$ .

```
upon event  $\langle beb, Broadcast \mid m \rangle$  do
  forall  $q \in \Pi$  do
    trigger  $\langle pl, Send \mid q, m \rangle$ ;

upon event  $\langle pl, Deliver \mid p, m \rangle$  do
  trigger  $\langle beb, Deliver \mid p, m \rangle$ ;
```

- When a messages is sent, the module sends along the id assigned to the upper layer which requested the send.
- The receive function registered to the runtime read the id contained in the header of the messages and call the callback related to this id that was registered by the upper layer.

```
328 #ifndef BUFFEREDSEND
329     if (_remove_buffer() == -1)
330         return -1;
331
332     MPI_Request** mpi_requests = (MPI_Request**) malloc(sizeof(MPI_Request*) * num_procs);
333     if (mpi_requests == NULL) {
334         sprintf(log_buffer, "BEB ERROR: Cannot allocate the mpi requests pointer table.");
335         runtime_log(ERROR_BEBROADCAST);
336         return -1;
337     }
338     if (_add_buffer(message, mpi_requests) == -1)
339         return -1;
340
341 #endif
342
343     for (current_process = 0; current_process < num_procs; current_process++) {
344         if (detected[current_process] == 0) {
345 #ifndef BUFFEREDSEND
346             MPI_Isend(message, size, MPI_INT, current_process, runtime_receive_handle, comm,
347                     mpi_requests[current_process]);
348 #endif

```

```
467     while (parent->next != NULL) {
468         buffer_data * bd = parent->next;
469         flag = 0;
470         for (i = 0; i < num_procs; i++) {
471             if (detected[i] == 0) {
472                 MPI_Test(bd->mpi_requests[i], &flag, &status);
473                 if (!flag || (status.MPI_SOURCE == MPI_ANY_SOURCE && status.MPI_TAG == MPI_ANY_TAG)) {
474                     flag = 0;

```

The function that handles send requests uses the asynchronous semantic (MPI\_Isend) instead of the buffered mode since it is hard to predict the size of the buffer to be attached. However this requires to manage the MPI\_Request struct for each send. Each send is recorded in a data structure. When this is visited, the module checks if the relates send is terminated, if it has it deletes the record.

# Reliable Broadcast

- RB1: Validity: If a correct process  $p$  broadcasts a message  $m$ , then  $p$  eventually delivers  $m$ .
- RB2: No duplication: BEB2
- RB3: No creation: BEB 1
- RB4: Agreement: If a message  $m$  is delivered by some correct process, then  $m$  is eventually delivered by every correct process.

```
379 int add_to_delivered(unsigned int id, int real_sender) {
380     // Delivered (removed ids)
381     if (isAllowedSerialNumber(proc_data_set[real_sender]->gc_threshold, id) == -1) {
382         return -1;
383     }
384
385     delivered * father = proc_data_set[real_sender]->delivered_set;
386
387     // Skips messages with id < threshold, until the list ends or the message with id = 0 is found
388     while (father->next != NULL && isLowerSerialNumber(proc_data_set[real_sender]->gc_threshold, father->next->id) == 0) {
389         father->next = father->next->next;
390     }
391
392     // Removes oldest contiguous messages from delivered set
393     while (father->next != NULL && isNextSerialNumber(proc_data_set[real_sender]->gc_threshold, father->next->id) == 0) {
394         delivered * d = father->next;
395
396         gettimeofday(&gettimeofdayArg, NULL);
397         incrementNextSerialNumber(proc_data_set[real_sender]->gc_threshold, father->next->id, gettimeofdayArg.tv_sec);
398
399         father->next = father->next->next;
400         free(d);
401     }
402 }
```

The Reliable Broadcast implements two functions:

- The send function uses the BEB layer to send a broadcast, the init has taken care to get the id used by BEB to identify messages for this module
- A function which is registered to the BEB layer is called when a message for this module has been received. This checks if the received messages has already been delivered through a mechanism that uses the second mean of managing Serial Number.

```
upon event { rb, Broadcast | m } do
    trigger { beb, Broadcast | [DATA, self, m] };

upon event { beb, Deliver | p, [DATA, s, m] } do
    if m ∉ delivered then
        delivered := delivered ∪ {m};
        trigger { rb, Deliver | s, m };
        trigger { beb, Broadcast | [DATA, s, m] };
```

```
243 int _eager_reliable_broadcast_send(int *message, int count, int handle) {
244     sprintf(log_buffer, "EAGER RELIABLE BROADCAST INFO: Registering message \"%d\\n\"");
245     runtime_log(DEBUG_EAGERRELIABLEBROADCAST);
246
247     // Copying the message, adding EAGER RELIABLE BROADCAST callback.
248     int *buffer = (int*) malloc(sizeof(int) * (count + 3));
249     buffer[0] = handle;
250     buffer[1] = serial_number++;
251     buffer[2] = my_rank;
252     copy_to_array(&buffer[3], message, count);
253     _best_effort_broadcast_send(buffer, count + 3, best_effort_broadcast_handle);
254     free(buffer);
255     return 0;
256 }
```

# Consensus

- C1: Termination: Every correct process eventually decides some value.
- C2: Validity: If a process decides  $v$ , then  $v$  was proposed by some process.
- C3: Integrity: No process decides twice.
- C4: Agreement: No two correct processes decide differently..

The Consensus X functions:

- The propose function gets an array of values, sorts it and delete duplicates before sending it through the BEB layer along with the id given by the calling layer and the id wrt the module.
- A function is registered to the BEB layer and when it is called reads the consensus id and either retrieves the instances data or creates a new instances with the id.
- When a propose or a crash event is received, a function is called which checks if the decide conditions have been reached, retrieves the latest proposed set, sends it as decided set through the BEB layer and triggers the decided event calling the callback registered by the layer which requested the propose.

The managing of the consensus id is done though the first mechanism explained earlier

```
upon event { P, Crash | p } do
    correct := correct \ {p};

upon event { c, Propose | v } do
    proposals[1] := proposals[1] ∪ {v};
    trigger { beb, Broadcast | [PROPOSAL, 1, proposals[1]] };

upon event { beb, Deliver | p, [PROPOSAL, r, ps] } do
    receivedfrom[r] := receivedfrom[r] ∪ {p};
    proposals[r] := proposals[r] ∪ ps;

upon correct ⊆ receivedfrom[round] ∧ decision = ⊥ do
    if receivedfrom[round] = receivedfrom[round - 1] then
        decision := min(proposals[round]);
        trigger { beb, Broadcast | [DECIDED, decision] };
        trigger { c, Decide | decision };
    else
        round := round + 1;
        trigger { beb, Broadcast | [PROPOSAL, round, proposals[round - 1]] };

upon event { beb, Deliver | p, [DECIDED, v] } such that p ∈ correct ∧ decision = ⊥ do
    decision := v;
    trigger { beb, Broadcast | [DECIDED, decision] };
    trigger { c, Decide | decision };
```

# Total Order Broadcast

The message id in the delivered data structure are managed as in the reliable module

- TOB1: Validity: If a correct process  $p$  broadcasts a message  $m$ , then  $p$  eventually delivers  $m$ .
- TOB2: No duplication: No message is delivered more than once.
- TOB3: No creation: If a process delivers a message  $m$  with sender  $s$ , then  $m$  was previously broadcast by process  $s$ .
- TOB4: Agreement: If a message  $m$  is delivered by some correct process, then  $m$  is eventually delivered by every correct process.
- TOB5: Total order: Let  $m_1$  and  $m_2$  be any two messages and suppose  $p$  and  $q$  are any two correct processes that deliver  $m_1$  and  $m_2$ . If  $p$  delivers  $m_1$  before  $m_2$ , then  $q$  delivers  $m_1$  before  $m_2$ .

The Total Order Broadcast module registers handlers both at the Reliable broadcast layer and at the Consensus layer.

- When the send function is called, the module sends the message through RB layer.
- When the function registered to the RB layer is called, the given message is stored in the unordered data structure.
- If unordered is not empty and no consensus instances are running, a consensus instances with the id (wrt TOB) of the messages in unordered is launched.
- When the consensus returns the decided set, the decided id are removed from unordered and added to delivered and then the deliver event is triggered, calling the callback registered by the calling layer.

```
upon event { tob, Broadcast | m } do
  trigger { rb, Broadcast | m };

upon event { rb, Deliver | p, m } do
  if m ∉ delivered then
    unordered := unordered ∪ {(p, m)};

upon unordered ≠ ∅ ∧ wait = FALSE do
  wait := TRUE;
  Initialize a new instance c.round of consensus;
  trigger { c.round, Propose | unordered };

upon event { c.r, Decide | decided } such that r = round do
  forall (s, m) ∈ sort(decided) do // by
    trigger { tob, Deliver | s, m };
  delivered := delivered ∪ decided;
  unordered := unordered \ decided;
  round := round + 1;
  wait := FALSE;
```

# Bibliography

- [1] Reliable and Secure Distributed Programming (Second Edition), Christian Cachin, Rachid Guerraoui, Luis Rodrigues
- [2] Parallel Programming in C with MPI and OpenMP, Michael J. Quinn
- [3] An Introduction to Parallel Programming, Peter Pacheco
- [4] MPI: The Complete Reference <http://www.netlib.org/utk/papers/mpi-book/mpi-book.html>
- [5] Setting up a Beowulf Cluster Using Open MPI on Linux, <http://techtinkering.com/2009/12/02/setting-up-a-beowulf-cluster-using-open-mpi-on-linux/>
- [6] C Programming language (Second Edition), Brian W. Kernighan, Dennis M. Ritchie
- [7] Open MPI: Open Source High Performance Computing, <http://www.open-mpi.org/>