# Exercise Session 8
# GM and VSC

## Problem 1

Show that P is the weakest failure detector for Group Membership. The failure detector D is weakest for solving some problem A (e.g., Consensus or NBAC) if D provides the smallest amount of information about failures that allows to solve A.

**Answer.** In order to show that P is the weakest failure detector for Group Membership, we need to show that:

- P can be used to implement Group Membership.

- Group Membership can be used to implement P.

The first direction stems directly from the Group Membership implementation in the class. For the second direction, we assume that all processes run Group Membership algorithm. Whenever a new view is installed, all processes that are freshly removed from the view are added to the suspected set. This approach satisfies both Strong Completeness and Strong Accuracy of P, directly from the corresponding properties of Group Membership.

## Problem 2

In this problem we will change the view-synchronous communication (VSC) abstraction in order to allow joins of new processes. Answer to the following questions:

1. Are the properties of VSC (as given in the class) suitable to accommodate the joins of new processes. Why / Why not?

2. Change the properties of VSC, so that they allow for implementations that support the joins of new processes. (Hint: focus on the properties of group membership)

3. Sketch the changes we need to perform on the Consensus-based (Algorithm II) implementation of VSC in order to support joins.

**Answer.**
*Solution 2.1*
No, the properties are not suitable for joins. The most obvious property is Local Monotonicity. Joins imply that the set of correct processes in a view can increase, and this would break the local monotonicity property. Furthermore, Completeness and Accuracy only refer to crashes, without imposing any conditions on the correctness of joins.

*Solution 2.2*
First, we need to add a $< Join|p >$ event to allow new processes to join the group. After a process emits such an event, we says that it requested to join. The VSC layer emits a $< JoinOk >$ event to the application when it has successfully joined a view. The application can start emitting broadcast requests after it receives the JoinOk event.

Group membership properties Let us first look at the four group membership properties. View Monotonicity. The monotonicity property of VSC (GM1) ensures that the number of processes in a view decreases over time. Since new processes can join, this needs to change: We conside three possibilities:

- Get rid of this property entirely.

- Require that views do not change for nothing: If a process installs views (j, N) and (j + 1, M), then $M \neq N$.

- Require that views do not oscillate (i.e., travel back in time): if a process p installs views (i, M) and (j, N) where j > i, q ∈ M, and q ∉ N, then for all k > j, if p installs (j,O), then q ∉ O.

With the second option, the new property ensures that consecutive views have different sets of processes, i.e., that the view cannot change if there is no change in the correct set of processes. Notice, however, that it is still possible for two views to have the same set of processes, e.g., if a processes joins and then crashes. It is also possible for a process to repeatedly be included and excluded from a view. With the third option, once a process is excluded from a view it can never come back.

**Uniform agreement.** The uniform agreement property of VSC (GM2) ensures that all processes install the same sequence of view. We will keep this property.

**Completeness.** If we choose the third version of monotonicity, then we can keep the completeness property of the group membership abstraction. If we choose one of the first two, we need to make some changes: Because the sequence of views is no longer monotonic, we need to strengthen a bit the completeness property of VSC (GM3): If a process p crashes, then there is $i \in N$ such that for all correct process q, if $j > i$ and q installs view (j, M), then $p \notin M$. To ensure that processes which want to join eventually join a view, we add the following completeness property: If a correct process p requests to join, then there is an integer i such that every correct process eventually installs view (i, M) such that $p \in M$.

**Accuracy.** If a process p installs views (i, M) and (i + 1, N) where $q \in M$ but $q \notin N$, then q has crashed. On top of those properties, we will also require that a process is included in a view only if it requested so.

**Validity.** If some process installs a view (i, M) and some process q is in M, then q previously requested to join or $q \in \Pi$.

*Broadcast properties*

Let us now look at the broadcast properties of VSC. Those are the same of for reliable broadcast (RB1,2,3,4). We have two options: either a process which joins needs to "catch-up" on all previously delivered messages, or a new process can just start with the messages of the first view in which it is included. If we choose the first option, then we can leave RB1,2,3,4 unchanged. If we choose the second option, then we need to relax Agreement (RB4) so that a process need to deliver only the messages sent in the view to which it participates: If message m is delivered by some correct process in view (i, M), then m is eventually delivered by all the process belonging to M. This way, if p $\notin$ M then p does not have to deliver m.

*View Synchrony*

Finally, we will keep the View Synchrony (VS) property as is.

*Solution 2.3*

The solution is described in Algorithm 1, 2 on the last two pages of this document. The changes to the regular algorithm are highlighted in red (note that we used the consensus algorithm that appears in the book it is similar in spirit to the version in the slides).

We add two new local variables to the algorithm: joined and crashed. The joined variable is a boolean flag that is set to true after the process successfully joins a view (is part of the view members). The joined flag differentiates the behavior of processes that are just attempting to join. The crashed variable is a local set that keeps track of crash events received from the failure detector. This set is useful in executions where a process p attempts to join and then crashes. If another correct process p2 sees the join attempt only after the crash notification, it needs to remember that it has already seen a crash of p and to disregard the join.

For most events, the only difference to the original algorithm is that we impose the condition joined = true for event handlers. Recall that such a conditional event handler means that the events are implicitly buffered until the condition becomes true (see the document describing the language used for module specification in Additional Material section on the course website). For example, the crash handler is now conditioned by joined = true. This means that any crash event received by the process while it is still joining will be buffered. The events will, however, be handled right after the process successfully joins a view.

The joining begins when the application emits a Join event (line 21). If the process has not joined yet and is not part of the initial set of processes in the view, the process broadcasts a JoinReq message to every other process. The JoinReq message can be seen as a dual of the crash event. It will be the job of the receiving correct processes (that are already view members) to handle the join and propose the addition of the joining process to a view.

Upon receiving a JoinReq message (line 23), processes will add the joining process to their correct set. Note that if the receiving process has already seen a crash of the joining process, the correct set will not be changed (p crashed will be $\phi$). Changing the correct set will trigger the handler at line 37 and initiate a view change. Processes that have seen the broadcast from the joining process will propose it in the new view member set. Since the joining process uses best-effort broadcast, correct processes will eventually receive the JoinReq broadcast message (if the joining process is also correct).

Another difference with the initial algorithm is that once a decision is taken in the consensus and a process moves to a new view, every process broadcasts the new view (both its member set and id). This broadcast is useful

for joining processes. If a joining process sees that it is part of a new view, it will initialize its view id, member set and correct set accordingly. Finally, the joining process sets the joined flag to true (meaning that it will handle all buffered events) and emits a JoinOk indication to the application.

**Algorithm 1** View synchrony with joins, first part

```
 1: Implements:
 2:     VSCJ (vscj)

 3: Uses:
 4:     UniformConsensus (ucons)
 5:     BestEffortBroadcast (beb)
 6:     PerfectFailureDetector (P)

 7: upon event ⟨vscj, Init⟩ do
 8:     (vid, M) := (0, Π)
 9:     correct := Π
10:     flushing := false; blocked := false; wait := false;
11:     pending := ∅; delivered := ∅; crashed := ∅
12:     forall m do ack[m] := ∅
13:     seen := [⊥]^N
14:     trigger ⟨vscj, View | (vid, M)⟩
15:     if self ∈ Π then
16:         joined := true
17:     else
18:         joined := false
19:     end if

20:     upon event ⟨vscj, Broadcast | m⟩ such that blocked = false ∧ joined = true do
21:         pending := pending ∪ (self, m)
22:         trigger ⟨beb, Broadcast | [DATA, vid, self, m]⟩

23:     upon event ⟨vscj, Deliver | p, [DATA, id, s, m]⟩ such that joined = true do
24:         if id = vid ∧ blocked = false then
25:             ack[m] := ack[m] ∪ {p}
26:             if (s, m) ∉ pending then
27:                 pending := pending ∪ (s, m)
28:                 trigger ⟨beb, Broadcast | [DATA, vid, s, m]⟩
29:             end if
30:         end if

31:     upon ∃(s, m) ∈ pending : M ⊆ ack[m] ∧ m ∉ delivered ∧ joined = true do
32:         delivered := delivered ∪ {m}
33:         trigger ⟨vscj, Deliver | s, m⟩

34:     upon event ⟨P, Crash | p⟩ such that joined = true do
35:         correct := correct \ {p}
36:         crashed := crashed ∪ {p}

37:     upon correct ≠ M ∧ flushing = false ∧ joined = true do
38:         flushing := true
39:         trigger ⟨vscj, Block⟩

40:     upon event ⟨vscj, BlockOk⟩ such that joined = true do
41:         blocked := true
42:         trigger ⟨beb, Broadcast | [PENDING, vid, pending]⟩

43:     upon event ⟨beb, Deliver | p, [PENDING, id, pd]⟩ such that id = vid ∧ joined = true do
44:         seen[p] := pd

45:     upon ∀p ∈ correct : seen[p] ≠ ⊥ ∧ wait = false do
46:         wait := true
47:         vid := vid + 1
48:         initialize a new instance uc.vid of uniform consensus
49:         trigger ⟨uc.vid, Propose | (correct, seen)⟩
```

**Algorithm 2** View synchrony with joins, second part

```
 1: upon event ⟨uc.id, Decide | M', S⟩ do
 2:     ∀p ∈ M' : S[p] ≠ ⊥ do
 3:         ∀(s, m) ∈ S[p] : m ∉ delivered do
 4:             delivered := delivered ∪ {m}
 5:             trigger ⟨vscj, Deliver | s, m⟩
 6:     flushing := false; blocked := false; wait := false
 7:     pending = ∅
 8:     ∀m do ack[m] := ∅
 9:     seen := [⊥]^N
10:     M := M'
11:     trigger ⟨vscj, View | (vid, M)⟩
12:     ∀p ∈ M do
13:         trigger ⟨beb, Broadcast | [NewView, vid, M]⟩

14:     upon event ⟨beb, Deliver | [NewView, vid', M']⟩ such that joined = false do
15:         if self ∈ M' then
16:             (vid, M) := (vid', M')
17:             correct := M
18:             joined := true
19:             trigger ⟨vscj, JoinOk⟩
20:         end if

21:     upon event ⟨vscj, Join | self⟩ such that joined = false do
22:         trigger ⟨beb, Broadcast | [JoinReq, self]⟩

23:     upon event ⟨beb, Deliver | [JoinReq, p]⟩ such that joined = true do
24:         correct := correct ∪ {p} \ crashed
```