

# ID2203 Tutorial 4 - Consensus

Cosmin Arad   Tallat M. Shafaat  
`icarad(@)kth.se`   `tallat(@)kth.se`

Handout: February 26, 2010

Due: March 3, 2010

## 1 Introduction

The goal of this tutorial is to understand and implement a Paxos uniform consensus algorithm, as well as an eventual leader election algorithm, for the partially-synchronous (fail-noisy) system model.

## 2 Assignment

In this assignment, you have to implement the following three algorithms/abstractions in Kompics:

- Paxos Uniform Consensus (PUC) - Algo. 5.7 in book.
- Abortable Consensus (AC) - Algo. 5.5 and 5.6 in book.
- Eventual Leader Detector (ELD) - Algo. 2.7 and 2.8 in book.

The algorithms have been re-written in this document with support for multiple consensus instances. This means that there can be multiple instances of consensus running at the same time. Each such instance is differentiated by an identifier called Paxos/consensus instance and denoted as *id* in the algorithms. Thus, when a node wants to propose a value, it should also provide the Paxos instance for which it wants to propose that value.

It is highly recommended that you read the afore-mentioned algorithms in the book before starting this assignment. Also note that PUC uses both AC and ELD.

In ELD, the function *select* can be implemented to return the *min* from the addresses passed to it. Thus, the lowest ranked alive process will become the leader.

Your implementation should provide, available via scenario and GUI, the following commands:

**Pi-j** This means that the node on which this command is executed should propose the value  $j$  for a consensus instance identified by  $i$ . For simplicity, both  $i$  and  $j$  should be integers.

**Dk** This means that the node  $n$  on which this command is executed should wait until it gets a decision to all previous (ongoing) proposals made by  $n$ , then wait for  $k$  milliseconds, and then process any further commands. If there is no previous (ongoing) proposal made by  $n$ , this command should simply wait for  $k$  milliseconds and then process any further commands.

**W** The node should sort all the decisions received thus far according to Paxos instance identifier and print them out. Please note that you should give enough delay before executing this command so that all the nodes have decided for all Paxos instances.

**Example:** P1-7:D100:P3-5:P4-9:D20000:W

The above command means that the node should propose value 7 in Paxos instance 1. After the node receives a decide for consensus instance 1, it should wait for 100 milliseconds and then propose value 5 in Paxos instance 3. Immediately following that (no waiting for a decision), the node should propose value 9 for Paxos instance 4. The node should then wait for Paxos instances 3 and 4 to finish. Next, the node should sleep for 20 seconds and then sort and print its ‘decides’ for all Paxos instances (even those that it didn’t initiate/propose in) to the console.

You can then compare the ‘decides’ that different nodes received.

### 3 Exercises

You have to send your source code and report answering the exercises given below.

**Exercise 1** Both Abortable Consensus (Algorithm 2 and Algorithm 3) and Paxos Uniform Consensus (Algorithm 4 and Algorithm 5) use a *seenIds* set to manage different concurrent instances of consensus. The *seenIds* set grows indefinitely. Explain in your written report, for each algorithm in part, how this set could be garbage collected.

**Exercise 2** Can PUC be used in a fail-recovery model? If so, under what condition/assumption? If not, why?

**Exercise 3** Construct a topology with 3 processes, use 1000ms latency between the nodes. Experiment with the following two scenarios:

- 1: D2200:P1-1
- 2: D2000:P1-2
- 3: D2000:P1-3

and

- 1: D2000:P1-1
- 2: D2200:P1-2
- 3: D2200:P1-3

What value is decided for Paxos instance 1 in these two scenarios? Give reasoning in your report why that particular value has been chosen.

**Exercise 4** Construct a topology with 2 processes. Assign link delays and initialize ELD period and increment in such a way that the two processes that initiate concurrent PUC proposals abort at least once in AC. Present the topology, the ELD parameter (Timedelay) and the operation sequences (ops) of the two processes and explain the execution.

**Exercise 5** In the original presentation of the Paxos algorithm, processes can have different roles: proposer, acceptor, and learner. For each event handler of AC and PUC say processes with what role are meant to execute the respective handler. In other words, what types of messages are sent and expected (and what events are triggered and expected) by each process role. Of course one process can act in more than one role. In fact, in our case, processes act in all roles, but interestingly, they need not to. Please refer to lecture 11 and the Paxos paper available on the course web page.

**Exercise 6** In this exercise, you will read a paper on Consensus: The Big Misunderstanding (available on the course web page) and you should answer the following questions.

**Question 6.1.** According to the paper, is atomic broadcast (total order broadcast) equivalent to consensus? What does that mean?

**Question 6.2.** According to the paper, are all notions of time impossible in an asynchronous model?

**Question 6.3.** Assume we can set the timeout value very high, e.g. 1 minute, and that would guarantee that the failure detectors will never suspect a correct node inaccurately. We have then circumvented the FLP impossibility. What practical implications does this have?

**Question 6.4.** According to the paper, is it possible to solve consensus in one step?

---

**Algorithm 1** Fail-noisy Eventual Leader Detector

---

**Implements:**

Eventual Leader Detector ( $\Omega$ ).

**Uses:**

PerfectPointToPointLinks (pp2p).

```
1: upon event  $\langle \text{Init} \rangle$  do
2:   leader := select( $\Pi$ );            $\triangleright$  deterministic function known in advance
3:   trigger  $\langle \text{trust} \mid \text{leader} \rangle$ ;            $\triangleright$  by all processes, e.g. min
4:   period := TimeDelay;
5:   for all  $p_i \in \Pi$  do
6:     trigger  $\langle \text{pp2pSend} \mid p_i, [\text{HEARTBEAT}] \rangle$ ;
7:   end for
8:   candidateset :=  $\emptyset$ ;
9:   startTimer(period);
10: end event

11: upon event  $\langle \text{Timeout} \rangle$  do
12:   newleader = select(candidateset);
13:   if (leader  $\neq$  newleader  $\wedge$  newleader  $\neq$  NIL) then
14:     period := period +  $\Delta$ ;
15:     leader := newleader;
16:     trigger  $\langle \text{trust} \mid \text{leader} \rangle$ ;
17:   end if
18:   for all  $p_i \in \Pi$  do
19:     trigger  $\langle \text{pp2pSend} \mid p_i, [\text{HEARTBEAT}] \rangle$ ;
20:   end for
21:   candidateset :=  $\emptyset$ ;
22:   startTimer(period);
23: end event

24: upon event  $\langle \text{pp2pDeliver} \mid p_j, [\text{HEARTBEAT}] \rangle$  do
25:   candidateset := candidateset  $\cup \{p_j\}$ ;
26: end event
```

---

---

**Algorithm 2** Abortable Consensus: Read Phase

---

**Implements:**

AbortableConsensus (ac).

**Uses:**

BestEffortBroadcast (beb);

PerfectPointToPointLinks (pp2p).

```
1: upon event  $\langle \text{Init} \rangle$  do
2:   seenIds :=  $\emptyset$ ;
3:   majority :=  $\lfloor N/2 \rfloor + 1$ ;
4: end event

5: procedure initInstance(id) is
6:   if (id  $\notin$  seenIds) then
7:     tempValue[id] := val[id] :=  $\perp$ ;
8:     wAcks[id] := rts[id] := wts[id] := 0;
9:     tstamp[id] := rank(self);
10:    readSet[id] :=  $\emptyset$ ;
11:    seenIds := seenIds  $\cup$  {id};
12:   end if
13: end procedure

14: upon event  $\langle \text{acPropose} \mid \text{id}, v \rangle$  do
15:   initInstance(id);
16:   tstamp[id] := tstamp[id] +  $N$ ;
17:   tempValue[id] :=  $v$ ;
18:   trigger  $\langle \text{bebBroadcast} \mid [\text{READ}, \text{id}, \text{tstamp}[\text{id}]] \rangle$ ;
19: end event

20: upon event  $\langle \text{bebDeliver} \mid p_j, [\text{READ}, \text{id}, \text{ts}] \rangle$  do
21:   initInstance(id);
22:   if (rts[id]  $\geq$  ts or wts[id]  $\geq$  ts) then
23:     trigger  $\langle \text{pp2pSend} \mid p_j, [\text{NACK}, \text{id}] \rangle$ ;
24:   else
25:     rts[id] := ts;
26:     trigger  $\langle \text{pp2pSend} \mid p_j, [\text{READACK}, \text{id}, \text{wts}[\text{id}], \text{val}[\text{id}], \text{ts}] \rangle$ ;
27:   end if
28: end event
```

---

---

**Algorithm 3** Abortable Consensus: Write Phase

---

```
1: upon event  $\langle pp2pDeliver \mid p_j, [NACK, id] \rangle$  do
2:   readSet[id] :=  $\emptyset$ ;
3:   wAcks[id] := 0;
4:   trigger  $\langle acReturn \mid id, \perp \rangle$ ;
5: end event

6: upon event  $\langle pp2pDeliver \mid p_j, [READACK, id, ts, v, sentts] \rangle$  do
7:   if (sentts = tstamp[id]) then
8:     readSet[id] := readSet[id]  $\cup$  {(ts, v)};
9:     if (|readSet[id]| = majority) then
10:      (ts, v) := highest(readSet[id]);  $\triangleright$  largest timestamp
11:      if (v  $\neq \perp$ ) then
12:        tempValue[id] := v;
13:      end if
14:      trigger  $\langle bebBroadcast \mid [WRITE, id, tstamp[id], tempValue[id]] \rangle$ ;
15:    end if
16:  end if
17: end event

18: upon event  $\langle bebDeliver \mid p_j, [WRITE, id, ts, v] \rangle$  do
19:   initInstance(id);
20:   if (rts[id] > ts or wts[id] > ts) then
21:     trigger  $\langle pp2pSend \mid p_j, [NACK, id] \rangle$ ;
22:   else
23:     val[id] := v;
24:     wts[id] := ts;
25:     trigger  $\langle pp2pSend \mid p_j, [WRITEACK, id, ts] \rangle$ ;
26:   end if
27: end event

28: upon event  $\langle pp2pDeliver \mid p_j, [WRITEACK, id, sentts] \rangle$  do
29:   if (sentts = tstamp[id]) then
30:     wAcks[id] := wAcks[id] + 1;
31:     if (wAcks[id] = majority) then
32:       readSet[id] :=  $\emptyset$ ; wAcks[id] := 0;
33:       trigger  $\langle acReturn \mid id, tempValue[id] \rangle$ ;
34:     end if
35:   end if
36: end event
```

---

---

**Algorithm 4** Paxos Uniform Consensus (part 1)

---

**Implements:**

UniformConsensus (uc).

**Uses:**

AbortableConsensus (ac);

BestEffortBroadcast (beb);

EventualLeaderDetector ( $\Omega$ ).

```
1: upon event  $\langle Init \rangle$  do
2:   seenIds :=  $\emptyset$ ;
3:   leader := false;
4: end event

5: procedure initInstance(id) is
6:   if (id  $\notin$  seenIds) then
7:     proposal[id] :=  $\perp$ ;
8:     proposed[id] := decided[id] := false;
9:     seenIds := seenIds  $\cup$  {id};
10:  end if
11: end procedure

12: upon event  $\langle trust \mid p_i \rangle$  do
13:   if ( $p_i = \text{self}$ ) then
14:     leader := true;
15:     for all id  $\in$  seenIds do
16:       tryPropose(id);
17:     end for
18:   else
19:     leader := false;
20:   end if
21: end event

22: upon event  $\langle ucPropose \mid id, v \rangle$  do
23:   initInstance(id);
24:   proposal[id] :=  $v$ ;
25:   tryPropose(id);
26: end event
```

---



---

**Algorithm 5** Paxos Uniform Consensus (part 2)

---

```
1: procedure tryPropose(id) is
2:   if (leader = true  $\wedge$  proposed[id] = false  $\wedge$  proposal[id]  $\neq \perp$ ) then
3:     proposed[id] := true;
4:     trigger  $\langle acPropose \mid id, proposal[id] \rangle$ ;
5:   end if
6: end procedure

7: upon event  $\langle acReturn \mid id, result \rangle$  do
8:   if (result  $\neq \perp$ ) then
9:     trigger  $\langle bebBroadcast \mid [DECIDED, id, result] \rangle$ ;
10:  else
11:    proposed[id] := false;
12:    tryPropose(id);
13:  end if
14: end event

15: upon event  $\langle bebDeliver \mid p_i, [DECIDED, id, v] \rangle$  do
16:   initInstance(id);
17:   if (decided[id] = false) then
18:     decided[id] := true;
19:     trigger  $\langle ucDecide \mid id, v \rangle$ ;
20:   end if
21: end event
```

---