

Caltech, CMS. CS/IDS 142: Lecture 8.2
 Distributed Fault-Tolerant Consensus: Paxos
 Mani Chandy
 18 November 2019

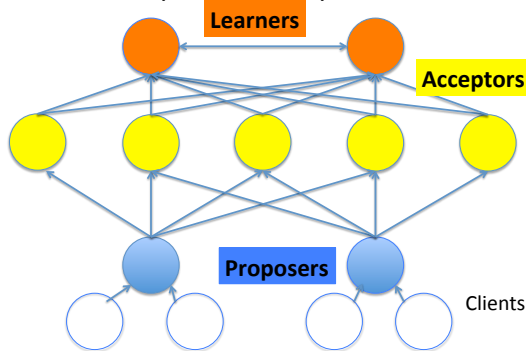
Take away concepts from the course:

- State transition systems.
- Always nothing bad happens.
- Progress. Variant function, Metric. The system never gets further away from its goal and eventually gets closer.
- Map global view to local agent state
- Data structures, e.g. graphs, that change with computation.

Consensus in Faulty Distributed Systems

- **Importance of consensus:** Multiple agents agree on a single sequence of events: banking, games, databases....
- **Modifications to usual distributed systems model:**
 - Agents may halt and restart from previously saved state.
 - Messages may be lost or duplicated.
- **Assumptions: No Byzantine operations**
 - Messages are not corrupted
 - Agents are not malicious

Clients, Proposers, Acceptors, Learners



Specification

Safety

1. Learners learn at most one value.

Associated with learner z is a variable $learned[z]$ where $learned[z]$ is either None (representing undefined) or some non-None value.

For $v \neq \text{None}$: $\text{stable}(learned[z] = v)$

2. Learners learn only proposed values

$(learned[z] \neq \text{None}) \Rightarrow learned[z]$ in set of proposals

3. All learners learn the same unchanging value

For all learners z, z' : $(learned[z]=\text{None})$ or $(learned[z']=\text{None})$ or $(learned[z]=learned[z'])$

Example: Group of students want to come to a consensus on an activity

- Some students in a group want to hike, others to swim, and others to dance.
1. Learners must learn a result proposed by a client. A learner must not learn that the result is "sleep" if the clients only propose "hike", "swim", "dance".
 2. All learners must learn the same result. All hike, or all swim, or all dance.
 3. After a learner learns a result that result remains unchanged. E.g. Learner won't switch from hike to swim.

Student is client, proposer, acceptor, and learner. The same agent plays all roles in many applications.

Progress Not Guaranteed!!

Not guaranteed: learners eventually learn a value.

For all learners z : $\text{eventually}(learned[z] \neq \text{None})$

Fischer, Lynch, Patterson (FLP) theorem says that **consensus cannot be achieved with a single faulty process**. Why?

We cannot prove progress; but we will discuss best effort algorithms.

Algorithm has two phases. A phase may be run multiple times.

- Phase 1:
 - Prepare(t) message from proposer to acceptor
 - promise(promise_t, accepted_t, value) reply from acceptor to proposer
- Phase 2:
 - Request(t, value) message from proposer to acceptor
 - Accept(t, value) message from acceptor to learner

Learner learns the proposal chosen by majority of acceptors.
 Different proposers never use the same t value.
 Break ties lexicographically

Algorithm for proposer P:

State: (P.t, P.value) Initially (-1, x)
 Start timer
 While not timed_out:

choose t greater than P.t and set P.t = t

PHASE 1

- send prepare(P.t) to all acceptors
- wait for promise(prepare_t, accept_t, value) replies from (at least) a majority of acceptors where prepare_t == P.t
- If value is not None for one or more of these promise messages then set P.value to the value in the promise message with the largest accept_t

PHASE 2

- send request(P.t, P.value) to all acceptors.
- Wait for accepted(t, value) replies from majority of acceptors where (t, value) == (P.t, P.value)

Algorithm for acceptor A:

State: (A.t, A.accepted_t, A.value) Initially (None, None, None)
 Start timer
 While not timed_out:

upon receiving prepare(t):
 if $t \geq A.t$:
 A.t = t
 reply with promise(t, A.accepted_t, A.value)

upon receiving request(t, value):
 If $t \geq A.t$:
 reply with accepted(t, value)

