

摘要

Basic Paxos 是分布式系统解决分布式协商问题的基本算法。很多人会忽视这个算法，认为它并不复杂且已经得到了很好的解决，但事实上算法本身在理解和解释上是存在一定困难的。本次讲座的内容按照Heidi Howard的博士论文<<Distributed consensus revised>>的内容，意图用作者自己的理解尽可能讲述Basic Paxos算法的过程，内在思想和一些安全性上的证明。

1 单值分布式协商问题

Single valued distributed consensus is the problem of deciding a single value $v \in V$ between a finite set of n participants, Non-triviality. The decided (commit) value must have been proposed by a participant.

Q1.1 为什么我们不能规定所有acceptor都选择同一个节点呢

Q1.2 为什么不能规定所有acceptor都按照同一顺序选择proposer的提案

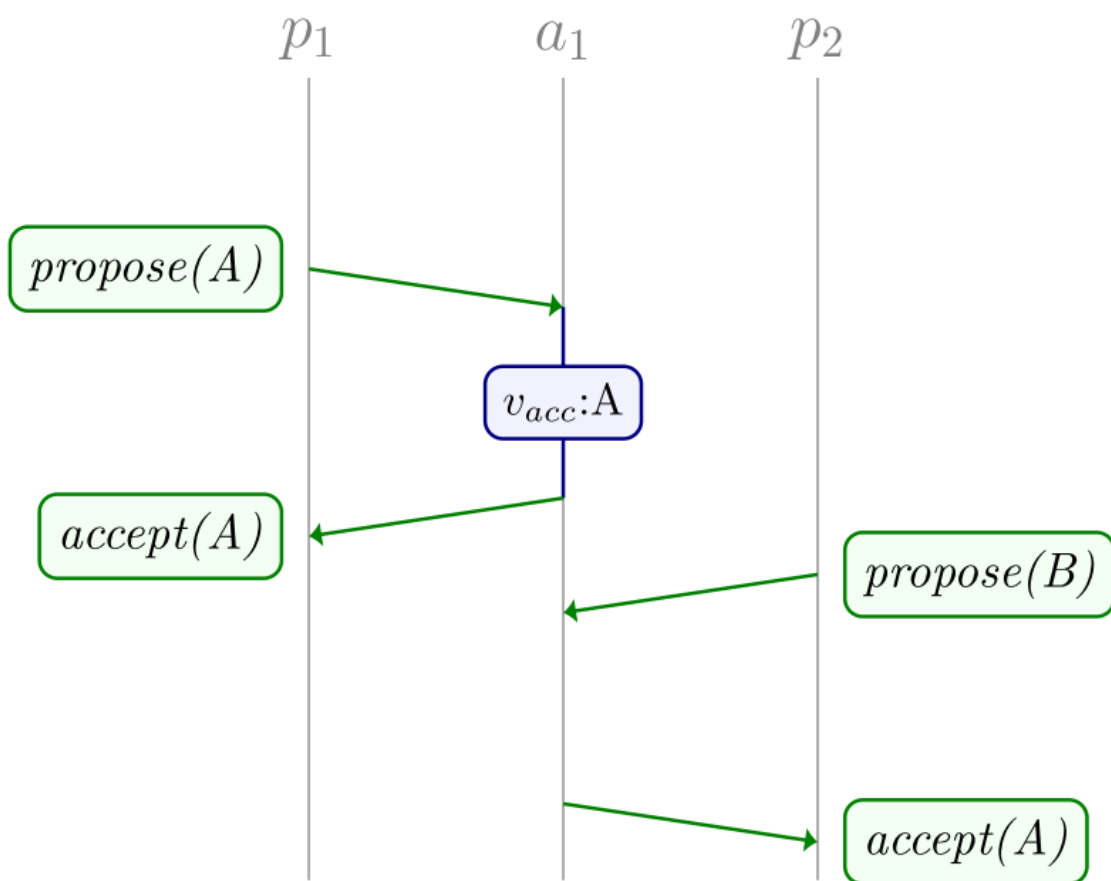
Safety. If a value has been decided, no other value will be decided.

Liveness. under certain liveness condition, it will eventually decide a value .

2 只有一个接受者的算法

idea: the acceptor determines the commit value as the first of all the proposals received.

example 2.1



algorithm 2.1

Algorithm 1: Proposer algorithm for SAA

state:

- γ : candidate value (configured, persistent)

```
1 send propose( $\gamma$ ) to acceptor
2 case accept( $v$ ) received from acceptor
    /* proposer learns that  $v$  was decided so return  $v$  */
3 return  $v$ 
```

Algorithm 2: Acceptor algorithm for SAA

state:

- v_{acc} : accepted value (persistent)

```
1 while true do
2     case propose( $v$ ) received from proposer
3         if  $v_{acc} = nil$  then
4              $v_{acc} \leftarrow v$ 
5         send accept( $v_{acc}$ ) to proposer
```

what are the difficulties if we have more than one acceptors ?

1. acceptors receive proposals in different sequence?

=> Q 2.1 why proposal order is required?

2. acceptors are not aware of the proposals identified by other acceptors

=> Q 2.2 why 2 phase algorithm is required?

3 Basic Paxos的过程

Phase 1: a proposer tries to learn the state of the system

Phase 2: a proposer tries to get a value to be accepted.

decided: if majority have accepted the value.

3. Q 3.1 what is the role of proposal id ?

after the value is decided?

before the value is decided?

algorithm 3.1

Classic Paxos Phase 1

1. A proposer chooses a unique epoch e and sends $prepare(e)$ to the acceptors.
2. Each acceptor stores the last promised epoch and last accepted proposal. When an acceptor receives $prepare(e)$, if e is the first epoch promised or if e is equal to or greater than the last epoch promised, then e is written to storage and the acceptor replies with $promise(e, f, v)$. (f, v) is the last accepted proposal (if present) where f is the epoch and v is the corresponding proposed value.
3. Once the proposer receives $promise(e, -, -)$ from the majority of acceptors, it proceeds to phase two. Promises may include a last accepted proposal which will be used by the next phase.
4. Otherwise if the proposer times out, it will retry with a greater epoch.

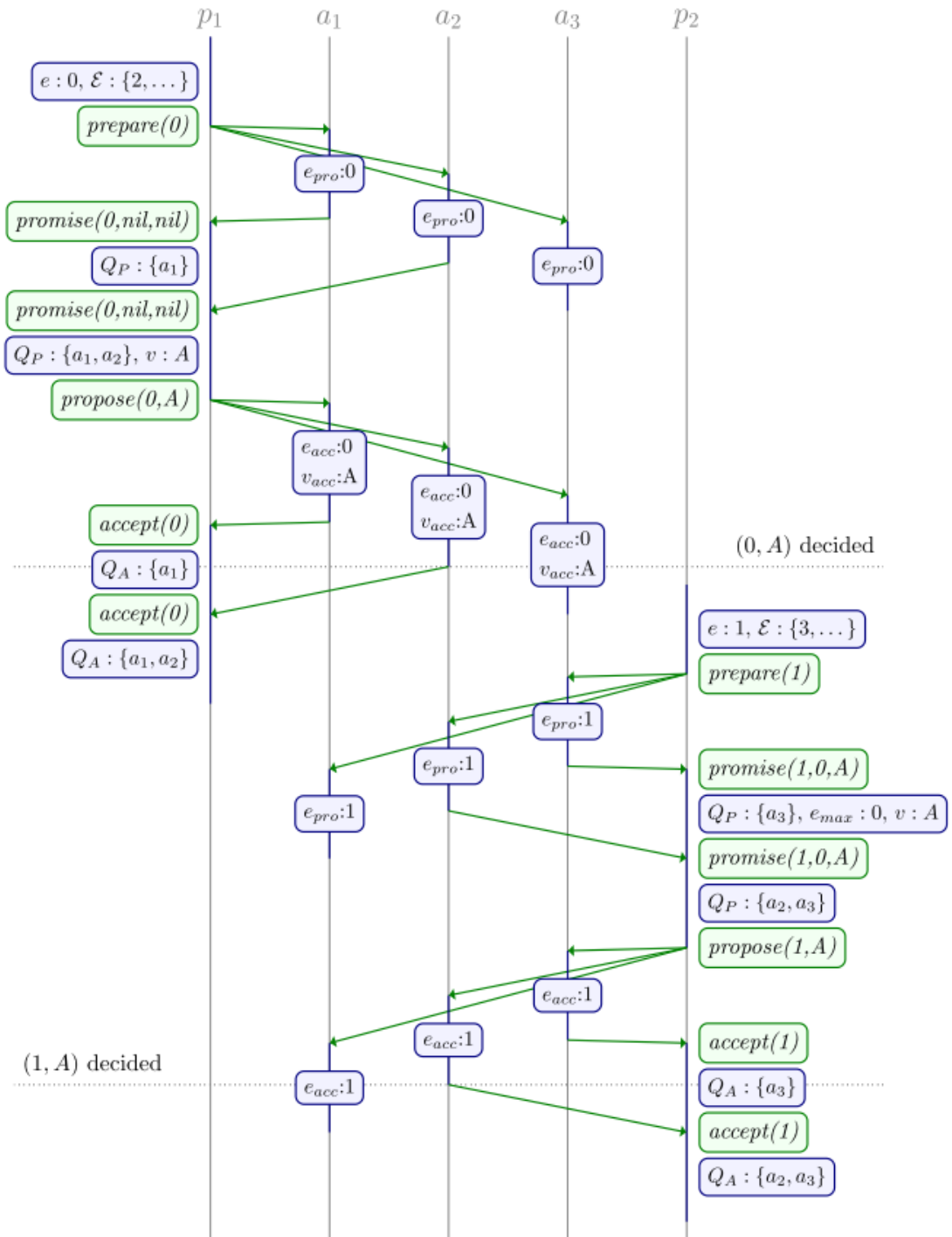
Classic Paxos Phase 2

1. The proposer must now select a value v using the following *value selection rules*:
 - i If no proposals were returned with promises in phase one, then the proposer will choose its candidate value γ .
 - ii If one proposal was returned, then its value is chosen.
 - iii If more than one proposal was returned then the proposer must choose the value associated with the greatest epoch.

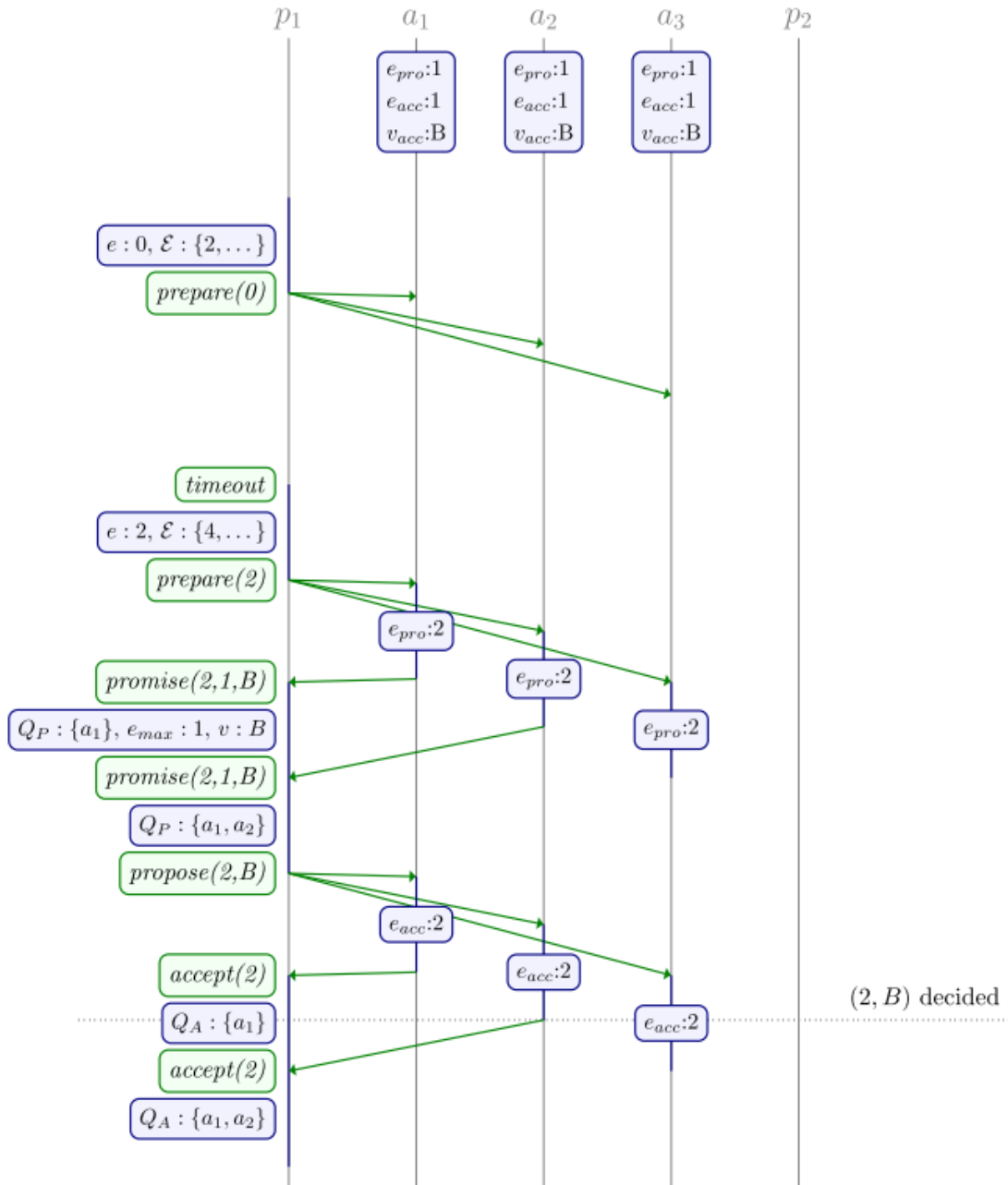
The proposer then sends $propose(e, v)$ to the acceptors.

2. Each acceptor receives a $propose(e, v)$. If e is the first epoch promised or if e is equal to or greater than the last promised epoch, then the promised epoch and accepted proposal is updated and the acceptor replies with $accept(e)$.
3. Once the proposer receives $accept(e)$ from the majority of acceptors, it learns that the value v is decided.
4. Otherwise if the proposer times out, it will retry phase 1 with a greater epoch.

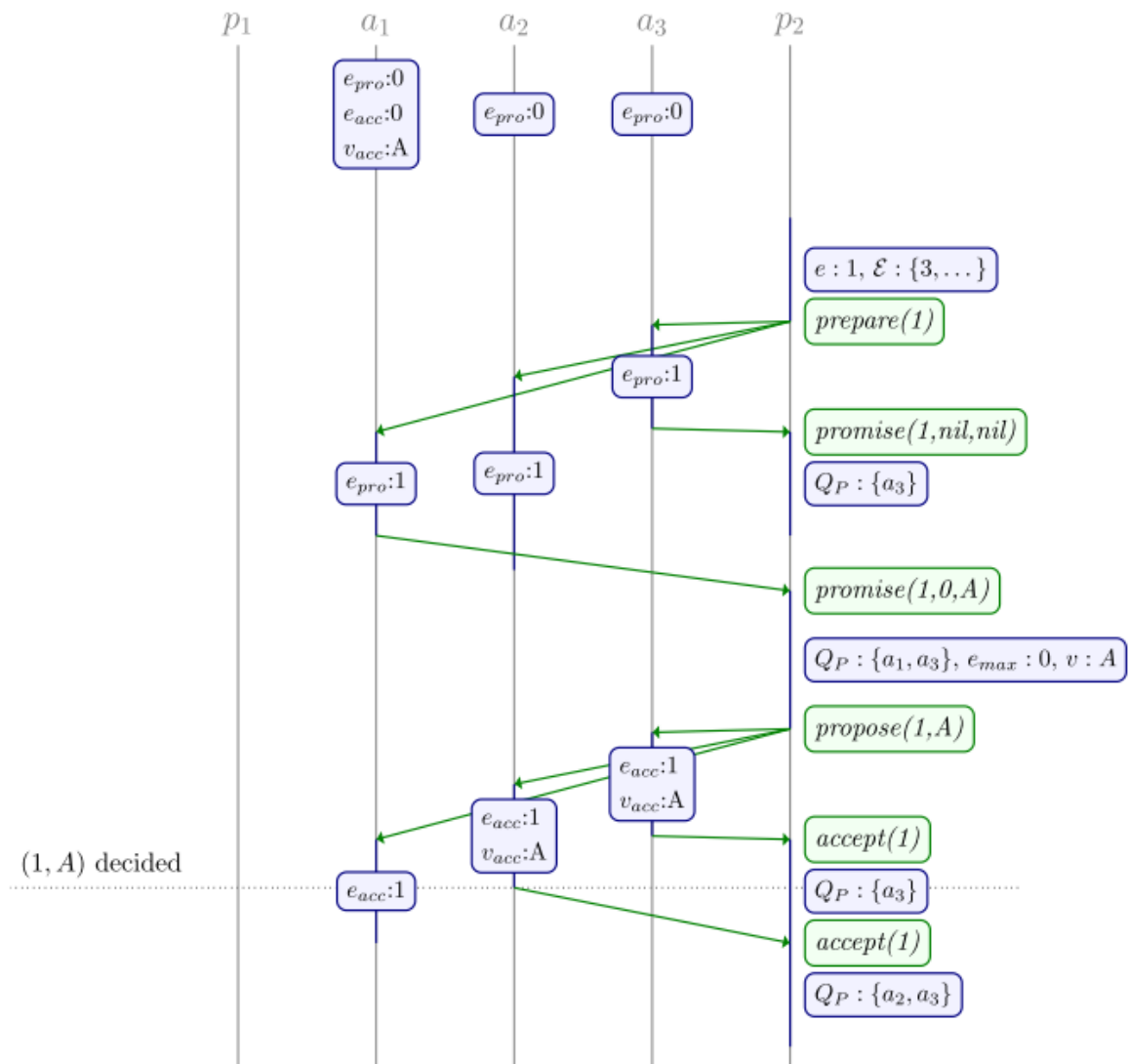
example 3.1



example 3.2



example 3.3



4. 安全性的理解

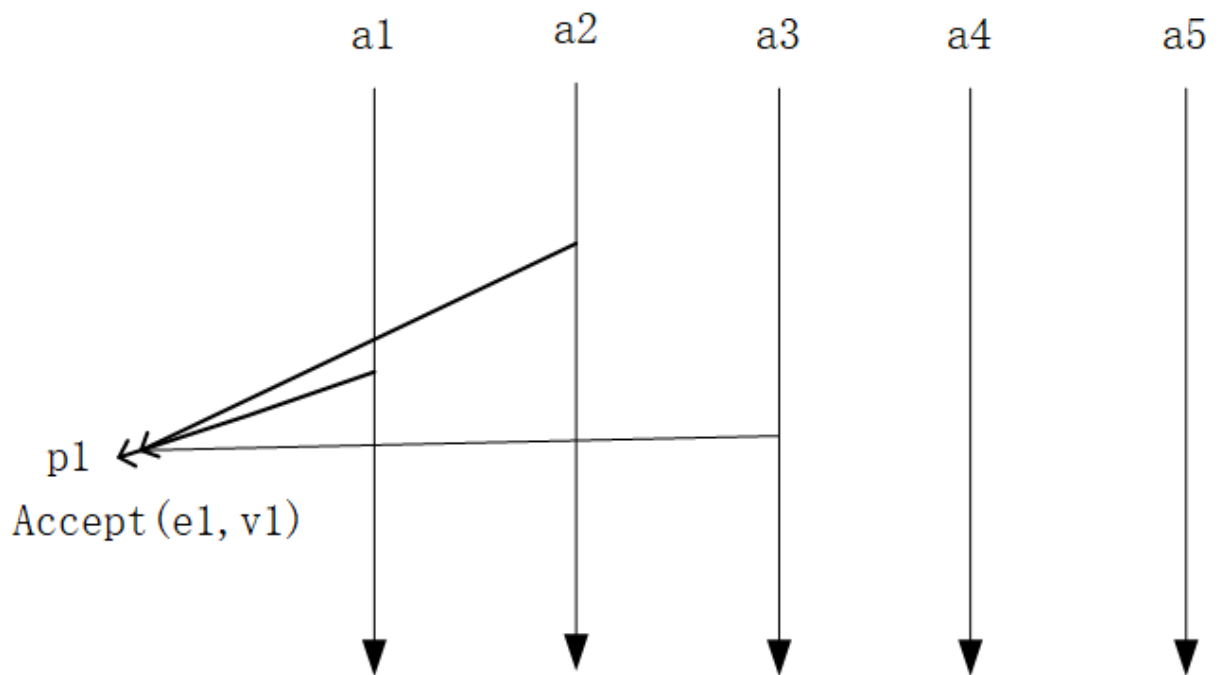
Q 4.1 how the algorithm guarantees that "If a value has been decided, no other value will be decided" ?

under the premise of:

至少一个存活的proposer 提出提案;

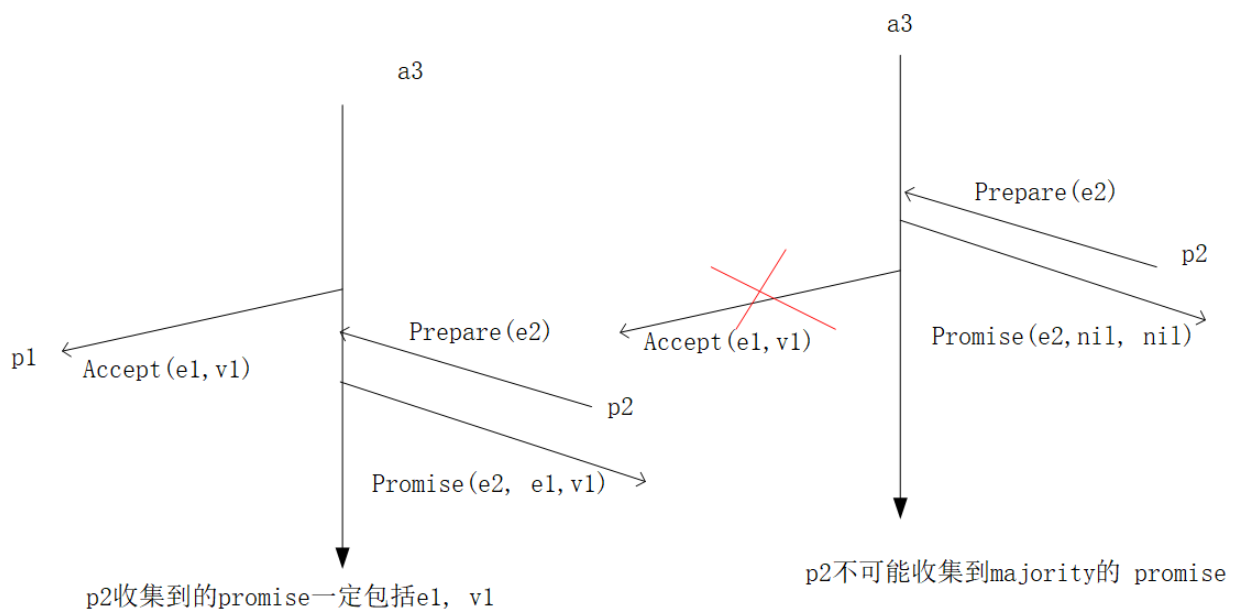
至少一半以上(majority) acceptor 存活;

假设如下图所示, 设 p_1 的proposal (e_1, v_1) 是第一个decided的proposal.



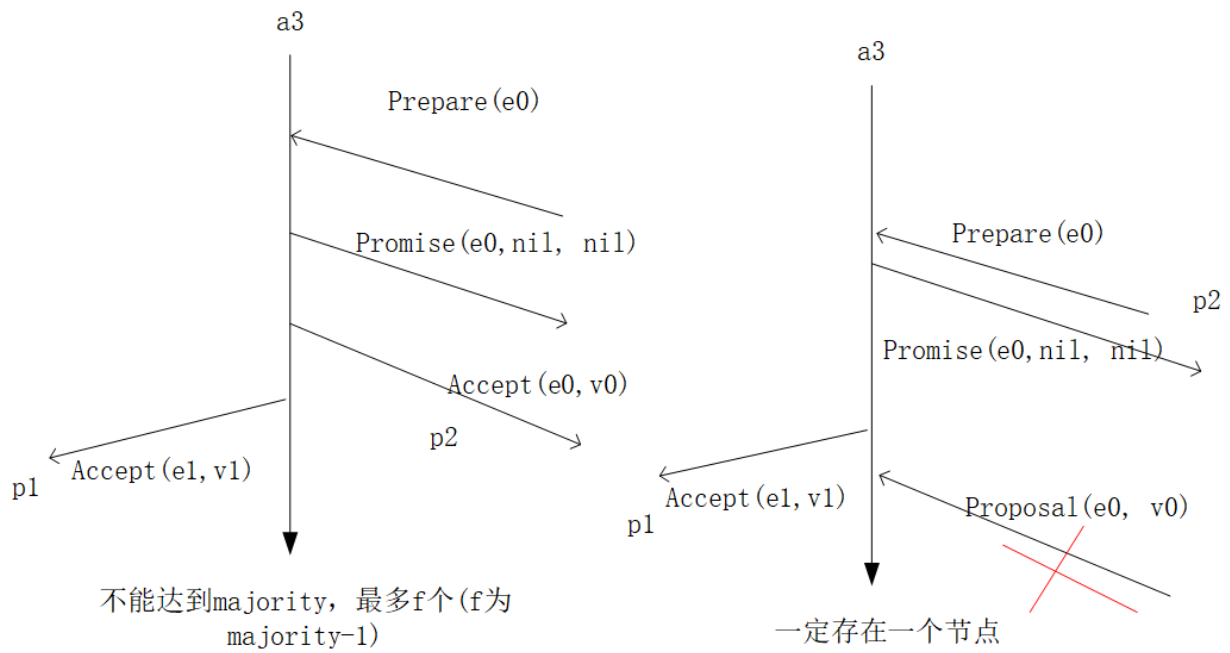
(1) 那么 $p1$ 的下一个proposal $p2$, 假设epoch 为 $e2$.

$\text{prepare}(e2)$ 的quorum和 $a1, a2, a3$ 必有一个交集, 假设为 $a3$



所以比 $e1$ 大的proposal不可能decide除 $v1$ 以外的一个value。

(2) 有没有可能比 $e1$ 小的proposal (比如 $e0$) 达成一致?



=Q4.1 再理解下epoch的作用?

如果 (e_1, v_1) 是第一个达成一致的proposal,

- (1) 那么比 e_1 小的proposal不可能达成一致, epoch的作用用来否决比 e_1 小的proposal;
- (2) 比 e_1 大的proposal 要么不可能达成一致, 要么达成一致的值为 v_1 。

5 Basic Paxos的一些优化

1. Distinguish Proposer.

参考文献

Heidi Howard. Distributed consensus revised

