

비동기적 분산 시스템하에서 선출 문제와 합의 문제의 관련성

(The Relationship between Election and Consensus in Asynchronous Distributed System)

박 성 훈 [†]

(Sung Hoon Park)

요 약 본 논문에서는 신뢰할 수 없는 고장추적 장치로 구성된 비동기적 분산 시스템 하에서 선출 (election) 문제와 합의(consensus) 문제의 관련성에 관하여 연구하고자 한다. 먼저 선출 문제는 합의 문제보다 더욱 어려운 문제임을 보인다. Chandra와 Toueg는 [8]에서 합의 문제는 비동기적 분산 시스템에서 신뢰할 수 없는 고장 추적 장치(unreliable failure detector)를 이용하여 해결 할 수 있음을 언급 하였다. 그러나, 합의 문제와는 대조적으로 선출 문제는 시스템 상에서 단 한 개의 노드가 죽은 경우에도 신뢰할 수 없는 고장 추적 장치를 이용하여 선출 문제를 해결할 수 없다. 이는 선출 문제는 합의 문제보다 더욱 어려운 문제임을 의미 한다. 보다 엄격하게 표현하자면, 선출 문제를 해결하는데 필요한 가장 약한 고장 추적 장치 (weakest failure detector)는 완전한 고장 추적 장치 (perfect failure detector) 이어야 하는 것으로, 이는 합의 문제를 해결하는데 필요한 가장 약한 고장 추적 장치보다 확실히 강한 것이다. 선출 문제가 합의 문제보다 어렵다는 것을 보이기 위해 축소(reduction) 프로토콜을 이용한다.

Abstract This paper discusses the relationship between the Election problem and the Consensus problem in asynchronous distributed systems with unreliable failure detectors. We first show that Election is harder to solve than Consensus. Chandra and Toueg have stated in [8] that Consensus is solvable in asynchronous systems with unreliable failure detectors. But, in contrast to the Consensus problem, the Election problem is impossible to solve with unreliable failure detectors even with a single crash failure. This means that the Election problem is harder than the Consensus problem. More precisely, the weakest failure detector that is needed to solve this problem is a Perfect Failure Detector, which is strictly stronger than the weakest failure detector that is needed to solve Consensus. We use a reduction protocol to show that the Election problem is harder than the Consensus problem.

1. Introduction

To elect a Leader (or Coordinator) in a distributed system, an agreement problem must be solved among a set of participating processes. This problem, called the Election problem, requires the participants to agree on only one leader in the system [1]. The problem has been widely studied

in the research community [2,3,4,5,6]. One reason for this wide interest is that many distributed protocols need an election protocol.

The Election problem is described as follows. At any time, there is at most one process that considers itself a leader and all other processes consider it as to be their only leader. If there is no leader, a leader is eventually elected.

Consensus and Election are similar problems in that they are both agreement problems. The so-called FLP impossibility result, which states that it is impossible to solve any non-trivial agreement in an asynchronous system even with a single

· 본 논문은 남서울대학교 2000년도 교내 연구비 지원으로 연구되었음.

· [†] 정 회 원 : 남서울대학교 컴퓨터학과 교수

spark@nsu.ac.kr

논문접수 : 2000년 3월 28일

심사완료 : 2000년 10월 26일

crash failure, applies to both problems [7]. The starting point of this paper is the fundamental result of Chandra and Toueg [8], which states that Consensus is solvable in asynchronous systems with unreliable failure detectors.

An interesting question is then whether the Election problem can also be solved in asynchronous systems with unreliable failure detectors. The answer to this question is No, and this is not surprising because the Election problem has been considered harder than Consensus [9]. However, in contrast to initial intuition, the reason Election is harder than Consensus is not its Liveness condition. The difficulty in solving Election is actually its Safety condition (all the nodes connected the system never disagree on the leader when the nodes are in a state of normal operation). This condition requires precise knowledge about failures which unreliable failure detectors cannot provide.

The rest of the paper is organized as follows. Section II describes motivations and the related works. In Section III we describe our system model. In Section IV we define Leader Election and show that it is harder than Consensus. Finally, Section V summarizes the main contributions of this paper and discusses related and future work.

2. Motivations and Related Works

In recent years, several paradigms have been identified to simplify the design of fault-tolerant distributed applications in a conventional static system. Election is among the most noticeable, particularly since it is closely related to group communication, which (among other uses) provides a powerful basis for implementing active replications.

It was shown in [7] that the Consensus problem cannot be solved in an asynchronous system if even a single crash failure can occur. The intuition behind this widely cited result is that in an asynchronous system, it is impossible for a process to distinguish between another process that has crashed and one that is merely very slow. The

consequences of this result have been enormous, because most real distributed systems today can be characterized as asynchronous, and Consensus is an important problem to be solved if the system is to tolerate failures.

As a result, the Consensus problem has frequently been used as a yardstick of computability in asynchronous fault-tolerant distributed systems. That means that if any problem is harder than Consensus, it also cannot be solved in asynchronous systems.

The asynchronous model of computation is especially popular in practice because unpredictable workloads are sources of asynchrony in many real systems. Therefore rendering any synchrony assumption is valid only probabilistically. Thus, the impossibility of achieving Consensus reveals a serious limitation of this model for fault-tolerant applications such as the Election problem. Because Consensus is such a fundamental problem, researchers have investigated various ways of circumventing the impossibility.

Actually, the main difficulty in solving such a problem in presence of process crashes lies in the detection of crashes. As a way of getting around the impossibility of Consensus, Chandra and Toueg extended the asynchronous model of computation with unreliable failure detectors and showed that the Consensus problem is solvable even with unreliable failure detectors [10].

If the Election problem is also solvable in asynchronous systems with unreliable failure detectors, it has an important consequence since the failure detection of a process is unreliable in real systems. To confirm whether Election is solvable in asynchronous systems with unreliable failure detectors, we compare Election with Consensus using a reduction protocol.

We are not the first to show that there are problems harder than Consensus. The first such result that we are aware of is [11] in which the authors show that Non-Blocking Atomic Commitment (NB-AC) cannot be implemented with the weakest failure detector that can implement

Consensus. This problem arises when transactions update data in a distributed system and the termination of transactions should be coordinated among all participants if data consistency is to be preserved even in the presence of failures [12]. It resembles the Election problem in that NB-AC is harder than Consensus.

To solve the NB-AC problem with an unreliable failure detector, they propose Non-Blocking Weak Atomic Commitment (NB-WAC) protocol and show that a failure detector weaker than a Perfect Failure Detector is strong enough to solve Non-Blocking Weak Atomic Commitment (NB-WAC). Hence, NB-AC appears to be harder than Consensus, but NB-WAC is easier than Election.

3. Model and Definitions

Our model of asynchronous computation with failure detection is the one described in [10]. In the following, we only recall some informal definitions and results that are needed in this paper.

3.1 Processes

We consider a distributed system composed of a finite set of processes $\Omega = \{p_1, p_2, \dots, p_n\}$ completely connected through a set of channels. Communication is by message passing, asynchronous and reliable. Processes fail by crashing; Byzantine failures are not considered.

Asynchrony means that there is no bound on communication delays or process relative speeds. A reliable channel ensures that a message, sent by a process p_i to a process p_j , is eventually received by p_j if p_i and p_j are correct (i.e. do not crash).

To simplify the presentation of the model, it is convenient to assume the existence of a discrete global clock. This is merely a fictional device inaccessible to processes. The range of clock ticks is the set of natural numbers. A history of a process $p_i \in \Omega$ is a sequence of events $h_i = e_i^0 \cdot e_i^1 \cdot e_i^2 \cdot \dots \cdot e_i^k$, where e_i^k denotes an events of process p_i occurred at time k . Histories of correct processes are infinite. If not infinite, the process history of p_i terminates with the event $crash_i^k$ (process p_i

crashes at time k). Processes can fail at any time, and we use f to denote the number of processes that may crash. We consider systems where at least one process is correct (i.e., $f < |\Omega|$).

A failure detector is a distributed oracle which gives hints on failed processes. We consider algorithms that use failure detectors. An algorithm defines a set of runs, and a run of algorithm A using a failure detector D is a tuple $R = \langle F, H, I, S, T \rangle$: I is an initial configuration of A ; S is an infinite sequence of events of A (made of process histories); T is a list of increasing time values indicating when each event in S has occurred; F is failure pattern that denotes the set $F(t)$ of processes that have crashed at any time t ; H is a failure detector history, which gives each process p and at any time t , a (possibly false) view $H(p, t)$ of the failure pattern: $H(p, t)$ denotes a set of processes, and $q \in H(p, t)$ means that process p suspects process q at time t .

3.2 Failure detector classes

Failure detectors are distributed oracles related to the detection of failures. A failure detector of a given class is a device that gives hints on a set of processes that it suspects to have crashed.

The Oracle notion has first been introduced as a language whose words can be recognized in one step from a particular state of a Turing machine [13,14]. The main characteristic of such oracles is to hide a sequence of computation steps in a single step (they may also hide an uncomputable function). They have been used to provide a hierarchy of problems. Hence the Oracle notation is related to the detection of failures. These oracles do not change the pattern of failures that affect the execution in which they are used. The main characteristic of such oracles is not related to the number of computation steps they hide, but to the guess they provide about failures.

Failure detectors are abstractly characterized by *completeness* and *accuracy* properties [10]. Completeness characterizes the degree to which crashed processes are permanently suspected by correct processes. Accuracy restricts the false suspicions

that a process can make.

Two completeness properties have been identified. *Strong Completeness*, i.e., there is a time after which every process that crashes is permanently suspected by every correct process, and *Weak Completeness*, i.e., there is a time after which every process that crashes is permanently suspected by some correct process.

Four accuracy properties have been identified. *Strong Accuracy*, i.e., no process is never suspected before it crashes. *Weak Accuracy*, i.e., some correct process is never suspected. *Eventual Strong Accuracy* (\diamond Strong), i.e., there is a time after which correct processes are not suspected by any correct process; and *Eventual Weak Accuracy* (\diamond Weak), i.e., there is a time after which some correct process is never suspected by any correct process. A failure detector class is a set of failure detectors characterized by the same completeness and the same accuracy properties (Figure 1).

For example, the failure detector class *P*, called *Perfect Failure Detector*, is the set of failure detectors characterized by Strong Completeness and Strong Accuracy. Failure detectors characterized by Strong Accuracy are reliable: no false suspicions are made. Otherwise, they are unreliable

Completeness	Accuracy			
	Strong	Weak	\diamond Strong	\diamond Weak
Strong	<i>P</i>	<i>S</i>	\diamond <i>P</i>	\diamond <i>S</i>
Weak	<i>Q</i>	<i>W</i>	\diamond <i>Q</i>	\diamond <i>W</i>

Fig. 1 Failure detector classes

For example, failure detectors of *S*, called Strong Failure Detector, are unreliable, whereas the failure detectors of *P* are *reliable*.

3.3 Reducibility and Transformation

The notation of *problem* reduction first has been introduced in the problem complexity theory [14], and in the formal language theory [13]. It has been also used in the distributed computing [15,16]. We consider the following definition of problem reduction.

An algorithm *A* solves a problem *B* if every run of *A* satisfies the specification of *B*. A problem *B* is said to be *solvable with* a class *C* if there is an algorithm which solves *B* using any failure detector of *C*. A problem *B*₁ is said to be reducible to a problem *B*₂ with class *C*, if any algorithm that solves *B*₂ with *C* can be transformed to solve *B*₁ with *C*. If *B*₁ is not reducible to *B*₂, we say that *B*₁ is *harder than B*₂.

A failure detector class *C*₁ is said to be *stronger than* a class *C*₂, (written *C*₁ ≥ *C*₂), if there is an algorithm which, using any failure detector of *C*₁, can emulate a failure detector of *C*₂. Hence if *C*₁ is stronger than *C*₂ and a problem *B* is solvable with *C*₂, then *B* is solvable with *C*₁. The following relations are obvious: *P* ≥ *Q*, *P* ≥ *S*, \diamond *P* ≥ \diamond *Q*, \diamond *P* ≥ \diamond *S*, *S* ≥ *W*, \diamond *S* ≥ \diamond *W*, *Q* ≥ *W*, and \diamond *Q* ≥ \diamond *W*. As it has been shown that any failure detector with *Weak Completeness* can be transformed into a failure detector with *Strong Completeness* [10], we also have the following relations: *Q* ≥ *P*, \diamond *Q* ≥ \diamond *P*, *W* ≥ *S* and \diamond *W* ≥ \diamond *S*. Classes *S* and \diamond *P* are incomparable.

3.4 Consensus

In the Consensus problem (or simply Consensus), every participant proposes an input value, and correct participant must eventually decide on some common output value [17]. Consensus is specified by the following conditions.

- *Agreement*: no two correct participant decide different values;
- *Uniform-Validity*: if a participant decides *v*, then *v* must have been proposed by some participant;
- *Termination*: every correct participant eventually decide.

Chandra and Toueg have stated that *Consensus* is solvable with \diamond *P* or *S* [10].

4. Election is harder than Consensus

In this section, we show that the Election problem is not solvable in asynchronous systems with unreliable failure detectors. This impossibility

result holds even with the assumption that at most one process may crash. Hence Election is harder than Consensus.

4.1 The Election Problem

Election is an important problem to solve for the construction of fault tolerant systems. It is closely related to the primary-backup approach (since choosing a primary replica is like electing a leader), an efficient form of passive replication. It is also closely related to group communication [18], which (among other uses) provides a powerful basis for implementing active replications.

The proof of the impossibility of Consensus in [7] assumes that it is impossible for a process to determine whether another process has crashed, or is just very slow. This assumption is widely cited as the reason for the impossibility result. There are other problems that cannot be solved in asynchronous systems with crash failures for the same intuitive reason that Consensus cannot be solved. Some of these problems can be solved with a weak failure detector; however, some cannot. In particular, the Election problem cannot be solved if a crashed process cannot be distinguished from a slow process.

The Election Problem is specified by the following two properties.

- *Safety*: All processes connected the system never disagree on a *leader* when the nodes are in a state of normal operation.

- *Liveness*: All processes should eventually progress to be in a state in which all processes connected to the system agree to the *only one* leader.

An *election protocol* is a protocol that generates runs that satisfies the Election specification.

4.2 Impossibility of solving Election Problem

Though $\diamond P$ or S are sufficient to solve Consensus, it is not sufficient to solve Election. Therefore the Election problem is strictly harder than the Consensus problem since even when assuming a single crash, unreliable failure detectors are not strong enough to solve election. In this section, we show that Strong Accuracy is necessary for solving Election, and it is sufficient

for solving Election.

Theorem 1 If $f > 0$, Election can not be solved with either $\diamond P$ or S .

PROOF. Consider a failure detector D of $\diamond P$. We assume for a contradiction that there exists a deterministic election protocol E that can be combined with a failure detector D such that $E + D$ is also an election protocol. Consider an algorithm A combined with $E + D$ which solves Election and a run $R = \langle F, H_D, I, S, T \rangle$ of A . We assume that only two processes P_i and P_j are correct and all messages from them is delayed until after t in R .

Consider that P_i is a leader at time (R, k) . At time (R, k_1) where $(k + t) > k_1 > k$, the process P_j falsely suspects other process P_i in some run. At time (R, k_2) where $k_2 > k_1$, P_j considers itself a leader by delaying the receipt of all messages sent by P_i until k_3 , where $(k + t) > k_3 > k_1$. Thus in (R, k_3) both P_i and P_j consider themselves the leader, violating the assumption that A is an election protocol.

But after a time t , all the processes except P_i and P_j are suspected. Hence there is a time after which every process that crashes is permanently suspected by every correct process. So H_D satisfies Strong Completeness. Consider Accuracy. After a time t , P_i and P_j are never suspected in H_D . Hence H_D satisfies Eventual Strong Accuracy. This is a contradiction.

Theorem 2 A weakest failure detector to solve Election is the Perfect Failure Detector.

PROOF: It is shown in [9] that a failure detector satisfying Strong Accuracy and Strong Completeness can be used to implement a Perfect Failure Detector. Strong Accuracy let processes never suspect a correct process: suspicions are never false. Every correct process always detects a leader failure only when the leader crashes using a Perfect Failure Detector. After an election is started, the problem of electing only one process as a leader is a kind of consensus problem; hence this problem is easily solved with a Strong Failure Detector that is less strong than Perfect Failure

detectors. That means that every correct process eventually gets into the state in which it considers only one process to be a leader. Therefore a Perfect Failure Detector is the weakest failure detector that is sufficient to solve Election.

5. Concluding Remarks

The importance of this paper is in extending the applicability field of the results, which Chandra and Toueg have studied on solving problems, into the Election problem in asynchronous system (with crash failures and reliable channels) augmented with unreliable failure detectors.

More specifically, what is the weakest failure detector in the asynchronous system? As an answer to this question, we showed that Perfect failure Detector P is the weakest failure detector to solve the Election problem in asynchronous systems. Though $\diamond P$ or S are sufficient to solve Consensus, we showed that they are not sufficient to solve Election. Therefore the Election problem is strictly harder than the Consensus problem even when assuming a single crash.

Determining that a problem Pb_1 is harder than a problem Pb_2 has a very important practical consequence, namely, the cost of solving Pb_1 cannot be less than that of solving Pb_2 . That means that the cost of solving Election cannot be less than that of solving Consensus.

The applicability of these results to problems other than Consensus has been discussed in [8, 17,18,19,20]. To our knowledge, it is however the first time that Election problems are discussed in asynchronous systems with unreliable failure detectors. We believe that there are problems harder than Election as well. One can define failure detectors that are stronger than a Perfect Failure Detector. For example, we can define a failure detector that is not only perfect but also guarantees that a failure of a process is detected only after all messages that it has sent have been received by the detecting process. This failure detector is required by some problems, including the non-

blocking version of the asynchronous Primary-Backup problem [12].

References

- [1] G. LeLann, "Distributed systems-towards a formal approach," in *Information Processing 77*, B. Gilchrist, Ed. North-Holland, 1977.
- [2] H.Garcia-Molian, "Elections in a distributed computing system," *IEEE Transactions on Computers*, vol. C-31, no. 1, pp. 49-59, Jan 1982.
- [3] H. Abu-Amara and J. Lokre, "Election in asynchronous complete networks with intermittent link failures." *IEEE Transactions on Computers*, vol. 43, no. 7, pp. 778-788, 1994.
- [4] H.M. Sayeed, M. Abu-Amara, and H. Abu-Avara, "Optimal asynchronous agreement and leader election algorithm for complete networks with byzantine faulty links.," *Distributed Computing*, vol. 9, no. 3, pp. 147-156, 1995.
- [5] J. Brunekreef, J.-P. Katoen, R. Koymans, and S. Mauw, "Design and analysis of dynamic leader election protocols in broadcast networks," *Distributed Computing*, vol. 9, no. 4, pp. 157-171, 1996.
- [6] G. Singh, "Leader election in the presence of link failures," *IEEE Transactions on Parallel and Distributed Systems*, vol. 7, no. 3, pp. 231-236, March 1996.
- [7] M. Fischer, N. Lynch, and M. Paterson. Impossibility of Distributed Consensus with One Faulty Process. *Journal of the ACM*, pp. 374-382. (32) 1985.
- [8] T. Chandra and S.Toueg. Unreliable failure detectors for reliable distributed systems. Technical Report, Department of computer Science, Cornell Univ., 1994.
- [9] D. Dolev and R Strong. A Simple Model For Agreement in Distributed Systems. In *Fault-Tolerant Distributed computing*, pp. 42-59. B. Simons and A. Spector ed, Springer Verlag (LNCS 448, 1987.
- [10] T. Chandra, V. Hadzilacos and S. Toueg. The Weakest Failure Detector for Solving Consensus. *Proceedings of the 11th ACM Symposium on Principles of Distributed Computing*, pp. 147-158. ACM press, 1992.
- [11] Rachid Guerraoui. Revisiting the relationship between non-blocking atomic commitment and consensus. In *Proceedings of the 10th International Workshop on Distributed Algorithms*, Springer Verlag (LNCS 857), 1996.
- [12] P.A.Bernstein, V. Hadzilacos, and N. Goodman.

- Concurrency Control and Recovery in Database Systems. Addison Wesley, 1987.
- [13] Hopcroft J.E. and Ullman J.D. *Introduction to Automata Theory, Languages and Computation*. Addison Wesley, Reading, Mass., 418 pages, 1979.
- [14] Garey M.R. and Johnson D.S. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. Freeman W.H & Co, New York, 340 pages, 1979.
- [15] Eddy Fromentin, Michel R RAY, Frederic TRONEL. On Classes of Problems in Asynchronous Distributed Systems. In *Proceedings of Distributed Computing Conference*. IEEE 10.4, June 1999.
- [16] Hadzilacos V. and Toueg S. Reliable Broadcast and Related Problems. In *Distributed Systems (Second Edition)*, ACM Press, New York, pp.97-145, 1993.
- [17] V. Hadzilacos. On the relationship between the atomic commitment and consensus problems. In *Fault-Tolerant Distributed Computing*, pp. 201-208. B. Simons and A. spector ed, Springer Verlag (LNCS 448), 1987.
- [18] A. Schiper and A. Sandoz. Primary Partition Virtually-Synchronous Communication harder than consensus. In *Proceedings of the 8th Workshop on Distributed Algorithms*, 1994.
- [19] L. Sabel and K. Marzullo. Election vs. Consensus in Asynchronous Systems. Technical Report TR95-1488, cornell Univ, 1995.
- [20] R. Guerraoui and A. Schiper. Transaction model vs Virtual Synchrony model: bridging the gap. In *Distributed Systems: From Theory to Practice*, pp. 121-132. K. Birman, F. Mattern and A. Schiper ed, Springer Verlag (LNCS 938), 1995.



박성훈

1982년 고려대학교 정경대학 통계학과 학사. 1992년 미국 Indiana Univ. (Bloomington) 전산학과 석사. 1994년 미국 Indiana Univ. (Bloomington) 전산학과 박사 수료. 1982년 ~ 1989년 두산그룹 기획조정실 전산실 팀장. 1994년 ~ 1996년 두산정보통신(주) 기술연구소 소장. 1996년 ~ 현재 남서울대학교 컴퓨터학과 조교수. 관심분야는 분산 알고리즘, 이동 컴퓨팅, 정형기법