Unreliable Failure Detectors for Reliable Distributed Systems

A different approach

- Augment the asynchronous model with an unreliable failure detector for crash failures
- Define failure detectors in terms of abstract properties, not specific implementations
- Identify classes of failure detectors that allow to solve Consensus

The Model

General

□ asynchronous system
 □ processes fail by crashing
 □ a failed process does not recover

Failure Detectors

- outputs set of processes that it currently suspects to have crashed
- □ the set may be different for different processes

Completeness

Strong Completeness Eventually every process that crashes is permanently suspected by every correct process

Weak Completeness Eventually every process that crashes is permanently suspected by some correct process

Accuracy

Strong Accuracy No correct process is ever suspected

Weak Accuracy Some correct process is never suspected

Accuracy

Strong Accuracy No correct process is ever suspected

Weak Accuracy Some correct process is never suspected

Eventual Strong Accuracy There is a time after which no correct process is ever suspected

Eventual Weak Accuracy There is a time after which some correct process is never suspected

Failure detectors

	Accuracy			
Completeness	Strong	Weak	Eventual strong	Eventual weak
Strong	Perfect P	$\operatorname{Strong} S$	$\Diamond P$	$\Diamond S$
Weak	Quasi Q	Weak W	$\Diamond Q$	$\Diamond W$



 $T_{\mathcal{D} \to \mathcal{D}'}$ transforms failure detector \mathcal{D} into failure detector \mathcal{D}'

If we can transform \mathcal{D} into \mathcal{D}' then we say that \mathcal{D} is stronger than \mathcal{D}' $(\mathcal{D} \ge \mathcal{D}')$ and that \mathcal{D}' is reducible to \mathcal{D}

If $\mathcal{D} \ge \mathcal{D}'$ and $\mathcal{D}' \ge \mathcal{D}$ then we say that \mathcal{D} and \mathcal{D}' are equivalent: $\mathcal{D} \equiv \mathcal{D}'$

Simplify, Simplify!

All weakly complete failure detectors are reducible to strongly complete failure detectors $P \ge Q, \quad S \ge W, \quad \Diamond P \ge \Diamond Q, \quad \Diamond S \ge \Diamond W$

Simplify, Simplify!

All weakly complete failure detectors are reducible to strongly complete failure detectors $P \ge Q, \quad S \ge W, \quad \Diamond P \ge \Diamond Q, \quad \Diamond S \ge \Diamond W$

All strongly complete failure detectors are
 reducible to weakly complete failure detectors (!)
 $Q \ge P$, W≥S, $\Diamond Q \ge \Diamond P$, $\Diamond W \ge \Diamond S$

Weakly and strongly complete failure detectors are equivalent!

From Weak Completeness to Strong Completeness

Every process p executes the following: $output_p := 0$ cobegin || Task 1: repeat forever ${p \text{ queries its local failure detector module } \mathcal{D}_p}$ $suspects_p := \mathcal{D}_p$ $send (p, suspects_p) \text{ to all}$ || Task 2: when receive($q, suspects_q$) from some q $output_p := (output_p \cup suspects_p) - {q}$ coend

The Theorems

Theorem 1 In an asynchronous system with W, consensus can be solved as long as $f \le n-1$

The Theorems

Theorem 1 In an asynchronous system with W, consensus can be solved as long as $f \le n-1$ Theorem 2 There is no f-resilient consensus protocol using $\Diamond P$ for $f \ge n/2$

The Theorems

Theorem 1 In an asynchronous system with W, consensus can be solved as long as $f \le n-1$ Theorem 2 There is no *f*-resilient consensus protocol using $\Diamond P$ for $f \ge n/2$

Theorem 3 In asynchronous systems in which processes can use $\Diamond W$, consensus can be solved as long as f < n/2

The Theorems

Theorem 1 In an asynchronous system with W, consensus can be solved as long as $f \le n-1$ Theorem 2 There is no *f*-resilient consensus protocol using $\Diamond P$ for $f \ge n/2$

Theorem 3 In asynchronous systems in which processes can use $\Diamond W$ consensus can be solved as long as f < n/2

Theorem 4 A failure detector can solve consensus only if it satisfies weak completeness and eventual weak accuracy-i.e. $\Diamond W$ is the weakest failure detector that can solve consensus.

Solving consensus using S

S: Strong Completeness, Weak Accuracy

- \Box at least some correct process c is never suspected
- O Each process p has its own failure detector
- Input values are chosen from the set {0,1}

Notation

We introduce the operators \oplus , \star , \otimes

They operate element-wise on vectors whose entries have values from the set $\{0, 1, \bot\}$

v * ⊥ = v	⊥ × v = v
v * v = ⊥	⊥∗⊥ = ⊥
v ⊕⊥ = v	⊥ ⊕v = v
$V \oplus V = V$	⊥⊕⊥=⊥
∨ ⊗⊥ = ⊥	⊥ ∞ ∨ = ⊥
$\mathbf{V} \otimes \mathbf{V} = \mathbf{V}$	⊥⊗⊥=⊥

Given two vectors A and B, we write A \leq B if A[i] $\neq \perp$ implies B[i] $\neq \perp$

Solving Consensus using any $\mathcal{D}\in S$

1: $V_p := (\bot, \dots, \bot, v_p, \bot, \dots, \bot)$	{p's estimate of the proposed values}		
2: $\Delta_p := (\bot, \dots, \bot, v_p, \bot, \dots, \bot)$			
3: {phase 1}	{asynchronous rounds $r_p, \ 1 \leq r_p \leq n-1$ }		
4: for $r_p := 1$ to $n-1$			
5: send (r_p, Δ_p, p) to all			
6: wait until [$\forall q$: received $(r_p, \Delta_q,$	$q) \; {\sf or} \; q \in \mathcal{D}_p {\sf J} $ {query the failure detector}		
7: $O_p := V_p$			
8: $V_p := V_p \oplus (\oplus_q received \Delta_q)$			
9: $\Delta_p := V_p \star O_p$	{value is only echoed the first time it is seen}		
10: {phase 2}			
11: send (r_p, V_p, p) to all			
2: wait until [$orall q$: received (r_p,V_q,q) or $q\in \mathcal{D}_p$]			
13: $V_p := \otimes_q received V_q$	{computes the "intersection", including V_p }		
14: {phase 3}			
15: decide on leftmost non- \perp coordinate of V_p			

A useful Lemma

 $\begin{aligned} &1: \ \ V_p := (\bot, ..., \bot, v_p, \bot, ..., \bot) & \{p$ is estimate of the proposed values} \\ &2: \ \ \Delta_p := (\bot, ..., \bot, v_p, \bot, ..., \bot) \end{aligned}$

- 3: {phase 1} {asynchronous rounds r_p , $l \le r_p \le n 1$ }
- 4: for rp := 1 to n-1
- 5: send (rp, ∆p,p) to all
- 6: wait until [$\forall q$: received (r_p, Δ_q, q) or $q \in \mathcal{D}_p$]
- 7: O_p := V_p
- 8: $V_p := V_p \oplus (\oplus_q \text{ received } \Delta_q)$
- 9: $\Delta_p := V_p \star O_p$ {value is only echoed first time it is seen} 10: {phase 2}
- 11: send (rp, Vp ,p) to all
- 12: wait until [$\forall q$: received (r_p , V_q , q) or $q \in \mathcal{D}_p$]
- 13: Vp := \otimes_q received Vq $\$ (computes the "intersection", including Vp) 14: (phase 3)
- 15: decide on leftmost non- \perp coordinate of V $_{p}$

Lemma 1 After phase 1 is complete, $V_c \leq V_p$ for all processes p that complete phase 1

A useful Lemma

- $\begin{array}{ll} 1: \ V_p \coloneqq (\bot, \, ..., \, \bot, \, v_p, \, \bot, \, ..., \, \bot) & \quad \ \left\{ p\text{'s estimate of the proposed values} \right\} \\ 2: \ \Delta_p \coloneqq (\bot, \, ..., \, \bot, \, v_p, \, \bot, \, ..., \, \bot) \end{array}$
- 3: {phase 1} {asynchronous rounds r_p , 1< $r_p \le n 1$ }
- 4: for r_p := 1 to n-1
- 5: **send** (r_p, ∆_p ,p) to all
- 6: wait until [$\forall q$: received (r_p, Δ_q, q) or $q \in \mathcal{D}_p$]
- 7: O_p := V_p
- $\mathsf{S:} \qquad \mathsf{V}_{\mathsf{p}} \coloneqq \mathsf{V}_{\mathsf{p}} \oplus (\oplus_{\mathsf{q}} \mathsf{received} \Delta_{\mathsf{q}})$
- 9: $\Delta_p := V_p \star O_p$ {value is only echoed first time it is seen} 10: {phase 2}
- 11: send (r_p, V_p, p) to all
- 12: wait until [$\forall q$: received (r_p , V_q , q) or $q \in \mathcal{D}_p$]
- 13: $V_p := \bigotimes_q \text{ received } V_q \text{ (computes the "intersection", including } V_p \}$ 14: {phase 3}
- 15: decide on leftmost non- \perp coordinate of V_p

Lemma 1 After phase 1 is complete, $V_c \leq V_p$ for all processes p that complete phase 1

Proof We show that

- $V_c[i] = v_i \land v_i \neq \bot \Rightarrow \forall p : V_p[i] = v_i$
- Let r be the first round when c sees v_i
- $\circ r \leq n-2$
- $\Box c$ will send to all Δ_c with v_i in round r
- □ By weak accuracy, all correct processes receive v_i by next round
- r = n 1
- \square v_i has been forwarded n-1 times: every other process has seen v_i

Two additional cool lemmas

- 8: $V_p := V_p \oplus (\oplus_q \text{ received } \Delta_q)$

9: $\Delta_p := V_p \star O_p$ {value is only echoed first time it is seen} 10: {Phase 2}

- 11: send (rp, Vp ,p) to all
- 12: wait until [$\forall q$: received (r_p , V_q , q) or $q \in D_p$]
- 13: Vp := \otimes_q received Vq $\{ \text{computes the "intersection", including Vp} \}$ 14: {Phase 3}
- 15: decide on leftmost non- \perp coordinate of V_p

Lemma 2 After Phase 2 is complete, $V_c = V_p$ for each p that completes phase 2

Proof

All processes that completed phase 2 have received V_c.
Since V_c is the smallest V vector, V_c[i] $\neq \bot \Rightarrow$ V_p[i] $\neq \bot \forall$ p

By the definition of ⊗
 $V_c[i] = ⊥ \Rightarrow V_p[i] = ⊥ \quad \forall p$

after phase 2

Lemma 3 $V_c \neq (\bot, \bot, \bot, \ldots, \bot)$

Solving consensus

1: $V_p := (\bot, ..., \bot, v_p, \bot, ..., \bot)$ {p's estimate of the proposed values} 2: $\Delta_p := (\bot, ..., \bot, v_p, \bot, ..., \bot)$ 3: {phase 1} for rp := 1 to n-1 send (r_{p}, Δ_{p}, p) to all wait until [$\forall q$: received (r_p , Δ_q , q) or $q \in \mathcal{D}_p$] $O_p := V_p$ $V_p := V_p \oplus (\oplus_q \text{ received } \Delta_q)$ 8: $\Delta_{\mathbf{p}} := \mathsf{V}_{\mathbf{p}} \star O_{\mathbf{p}} \qquad \{ \text{value is only echoed first time it is seen} \}$ 10: {phase 2} 11: send (rp, Vp,p) to all 12: wait until [$\forall q$: received (r_p , V_q , q) or $q \in \mathcal{D}_p$] 13: $V_p := \bigotimes_q received V_q$ {computes the "intersection", including 14: {phase 3} 15: decide on leftmost non- \perp coordinate of V_p

Theorem The protocol to the left satisfies Validity, Agreement, and Termination

Proof

Left as an exercise

A lower bound - I

Theorem Consensus with $\Diamond P$ requires f < n/2

A lower bound - I

Theorem Consensus with $\Diamond P$ requires f < n/2

Proof

- o Suppose n is even, and a protocol exists that solves consensus when $f\!=\!n/2$
- O Divide the set of processes in two sets of size $n/2, P_1$ and P_2







The case of the Rotating Coordinator

Solving consensus with $\Diamond W$ (actually, $\Diamond S$)

- Asynchronous rounds
- ${\ensuremath{\mathfrak{G}}}$ Each round has a coordinator c
- $\bigcirc c_{id} = (r \mod n) + 1$
- ${\ensuremath{\textcircled{}}}$ Each process p has an opinion $v_p \!\in\! \{0,1\}$ (with a time of adoption t_p)
- Ocoordinator collects opinions to form a suggestion
- If they believe c to be correct, processes adopt its suggestion and make it their own opinion
- A suggestion adopted by a majority of processes is "locked"

One round, four phases

Phase 1

Each process, including c , sends its opinion timestamped r to c

One round, four phases

Phase 1

Each process, including c, sends its opinion timestamped r to c

Phase 2

- c waits for first $\lceil n/2 + 1 \rceil$ opinions with timestamp r
- \boldsymbol{c} selects $\boldsymbol{v}\text{,}$ one of the most recently adopted opinions
- \boldsymbol{v} becomes \boldsymbol{c} 's suggestion for round \boldsymbol{r}
- \boldsymbol{c} sends its suggestion to all

One round, four phases

Phase 1

Each process, including c, sends its opinion timestamped r to c

Phase 2

- c waits for first $\lceil n/2 + 1 \rceil$ opinions with timestamp r
- c selects v, one of the most recently adopted opinions
- v becomes c 's suggestion for round r
- c sends its suggestion to all

Phase 3

Each p waits for a suggestion, or for failure detector to signal c is faulty If p receives a suggestion, p adopts it as its new opinion and ACKs to c Otherwise, p NACKs to c

One round, four phases

Phase 1

Each process, including c, sends its opinion timestamped r to c

Phase 2

- c waits for first $\lceil n/2+1 \rceil$ opinions with timestamp r
- c selects v, one of the most recently adopted opinions
- v becomes c's suggestion for round r
- c sends its suggestion to all

Phase 3

Each p waits for a suggestion, or for failure detector to signal c is faulty If p receives a suggestion, p adopts it as its new opinion and ACKs to cOtherwise, p NACKs to c

Phase 4

c waits for first $\lceil n/2+1 \rceil$ responses if all ACKs, then c decides on v and sends DECIDE to all if p receives DECIDE, then p decides on v

Consensus using $\Diamond S$

 $v_p :=$ input bit; r := 0; $t_p := 0$; $state_p :=$ undecided while p undecided do $c = (r \mod n + 1)$ {phase 1: all processes send opinion to current coordinator} p sends (p, r, v_p, t_p) to c {phase 2: current coordinator gather a majority of opinions} c waits for first $\lceil n/2+1 \rceil$ opinions (q, r, v_q, t_q) c selects among them the value v_a with the largest t_a c sends (c, r, v_q) to all {phase 3: all processes wait for new suggestions from the current coordinator} p waits until suggestion (c, r, v) arrives or $c \in \Diamond S_p$ if suggestion is received then $\{v_p := v; t_p := r; p \text{ sends } (r, ACK) \text{ to } c\}$ else p sends (r, NACK) to c(phase 4: coordinator waits for majority of replies. If majority adopted the coordinator's suggestion, then coordinator sends request to decide) c waits for first $\lceil n/2+1 \rceil$ (r, ACK) or (r, NACK) if c receives $\lceil n/2+1 \rceil$ ACKs, then c sends (r, DECIDE, v) to all

when p delivers (r, DECIDE, v) then $\{p \text{ decides } v ; state_p := \text{ decided}\}$

$\Diamond S$ Consensus as Paxos

- All processes are acceptors
- In round r, node (r mod n)+1 serves both as a distinguished proposer and as a distinguished learner
- The round structure guarantees a unique proposal number
- The value that a proposer proposes when no value is chosen is not determined