

Failure Detectors

Poster session: Next week, here

- No need to print poster: Just show it on your laptop, etc.
- Ideally: 1 slide/screen (Will NOT look at more than 2)
 - What were you doing
 - What was the best/coolest thing you found.
 - Things You WANT YOUR CLASSMATES TO KNOW
- MAIN AUDIENCE: OTHERS IN THE CLASS!

Where we are

RSM → Consensus/Agreement

↳ Validity, agreement, termination

Fail-Stop

Asynch

X

BFT (w/auth.)

X

Partially Synch

✓ if $n > 2f$

✓ if $n > 3f + 1$

Synch

✓ if $n > f + 1$

✓ if $n > f + 1$

Why gap? Hand-wavy claim

↳ Partial synchrony allows distinguishing

b/w failure & delay

↳ When?

→ For how many nodes?

→ ...

Can we more precisely define what we need?

Why - Theory: Better understand the protocol

Practice: Maintaining/establishing assumptions come at a cost

↳ $n > 2f$

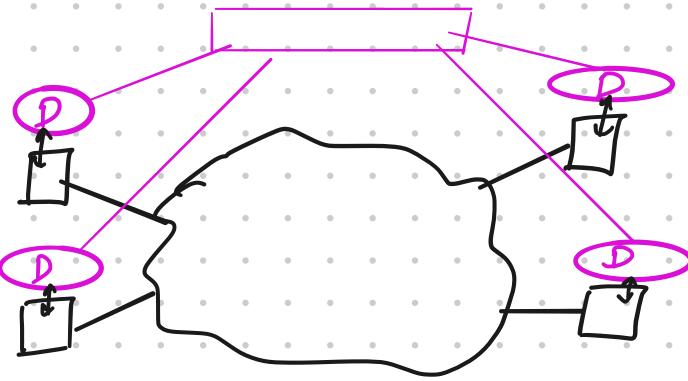
→ Partial synchrony

→ Bounded message latency

..

Useful to know what is required.

Failure Detectors



Modeling Notes

- F_0 Time → Output
 (Z^t) can be anything

In this paper

Output = 2^π : set of processes
suspected of failing

Other Work

Output = {Green, Red} [Green ⇒ No process
is suspected of
having failed]

= 2^π : Quorum

= 2^π / Green / Red (sufficient for NBAC)

Can we more precisely define what we need?



① Show consensus (or
other problem) can be



solved if each processes had access to F , assuming some failure model.

Chandra and Toueg :

Strong (Σ) or

$\Diamond W(\Sigma)$ w/ $2f+1$ processes suffice

② Show that if $\exists F'$ fodo to solve consensus then

one can construct F from F'

Chandra, Hadzilacos & Toueg (CHT) :

$\Diamond W(\Sigma)$ is weakest F.D. for

Solving consensus

[Roughly: Use F' to run a consensus protocol that tracks participants]

$\Diamond W \rightarrow \exists p \in T$ that everyone agrees has not failed

↳ Leader election?

Chandra & Toueg's Failure Detectors

F : Time \rightarrow Set of suspect processes

Accuracy: "What correct processes can be suspected"

- Strong: No correct process is ever suspected

- Weak: \exists correct process p that is never suspected.

- \Diamond Strong: Eventually no correct process is suspected

$\exists t \text{ s.t. } t' \geq t \Rightarrow F(t')$ only contains failed processes

- \Diamond Weak: \exists correct process p that is eventually not suspected

$\exists p \in T, t: t' \geq t \Rightarrow p \notin F(t')$

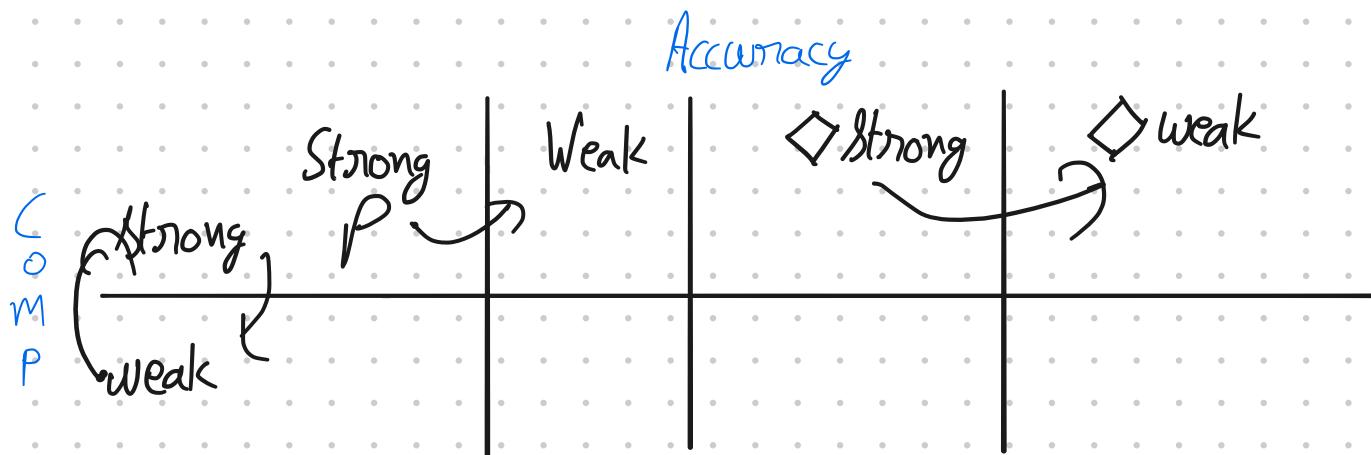
Completeness: "What faulty processes are suspected"

- Strong: Eventually, all failed processes are suspected by all processes

- Weak: Eventually, some correct process suspects each failed process.

Combine these two dimensions to come up with

8 F.D.s



Eight is too many objects.

Observe for

- Comp & accuracy

Strong \Rightarrow weak

- Accuracy

\diamondsuit Strong \Rightarrow \diamondsuit weak

Can build F.D. with strong completeness given
F.D. with weak completeness.

① Completeness is eventual

→ Can exchange messages and
get FD output from other
correct processes

Idea Strong complete output

[all correct processes eventually
suspect all failed processes]

||

U_{correct} weak complete output

P [a correct process eventually
suspects any failed process]

② Tricky: preserving accuracy.

a) Strong \leftarrow No correct process is
suspected

\Rightarrow No processes FD output contains
correct processes

\Rightarrow Union does not contain correct p

b) Weak \leftarrow I correct p that is never suspected

[same reasoning as above]

c) \diamondsuit strong \wedge \diamondsuit weak: Tricky

" " \Rightarrow 1. might contain correct

Why? FD output might contain incorrect processes (or chosen correct process p) initially

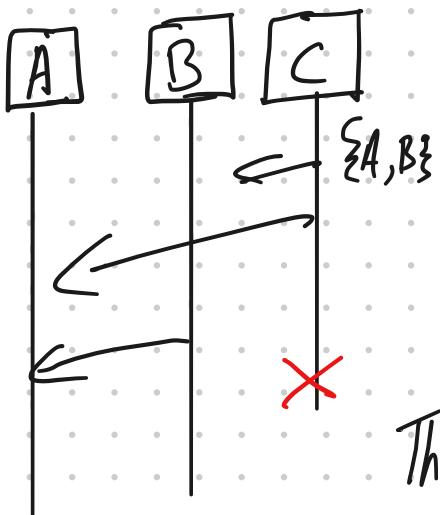
Need to get rid of it eventually.

How?

Idea 1

while true:

send ($\langle p, \text{fd output} \rangle$)



on receive $\langle p, f \rangle$

$out[p] \leftarrow f$

fd output = $\cup out[p]$

This doesn't work. Why?

Idea 2:

Every process p executes the following:

$output_p \leftarrow \emptyset$

cobegin

|| Task 1: repeat forever

{ p queries its local failure detector module D_p }
 $suspects_p \leftarrow D_p$
send $(p, suspects_p)$ to all

} Send FD output

|| Task 2: when receive $(q, suspects_q)$ for some q

$output_p \leftarrow (output_p \cup suspects_q) - \{q\}$

coend

Why does this work?

C
O
M
P

Accuracy

Strong	Weak	Strong	Weak
Strong	P	S	$\diamond P$
Weak	Q	W	$\diamond Q$

Can focus on one or the other.

Two results

- ① $S \text{ (or } W)$ → Consensus with $f+1$ nodes
- ② $\diamond S \text{ (or } \diamond W)$ → Consensus with $2f+1$ nodes.

Going to start with ② ← Similar to what we have seen before.

Core idea

- ① Nodes take it in turn to be leader

↳ When a node becomes leader

it tries to replicate a

value (proposal) to all connected nodes.

→ Leader counts how many times its value has been replicated \wedge decides (commits) when the value is sufficiently replicated
↳ Informs nodes about decisions

- ② Other processes move onto another leader if F.D. suspects current leader

Timing might lead to a scenario where current leader is suspected after value is sufficiently replicated (\wedge potentially committed)

↳ Must ensure that a sufficiently replicated value is used for all future proposals.

↳ Compare to Paxos

P2. If a proposal with value v is chosen, then every higher-numbered proposal that is chosen has value v .

Use similar idea \leftarrow Quorum Intersection

to decide proposal.

Every process p executes the following:

```

procedure propose( $v_p$ )
     $estimate_p \leftarrow v_p$                                 { $estimate_p$  is  $p$ 's estimate of the decision value}
     $state_p \leftarrow \text{undecided}$ 
     $r_p \leftarrow 0$                                      { $r_p$  is  $p$ 's current round number}
     $ts_p \leftarrow 0$                                     { $ts_p$  is the last round in which  $p$  updated  $estimate_p$ , initially 0}

```

{Rotate through coordinators until decision is reached}

while $state_p = \text{undecided}$

```

     $r_p \leftarrow r_p + 1$ 
     $c_p \leftarrow (r_p \bmod n) + 1$                          { $c_p$  is the current coordinator}

```

Phase 1: {All processes p send $estimate_p$ to the current coordinator}
 send $(p, r_p, estimate_p, ts_p)$ to c_p

Phase 2: {The current coordinator gathers $\lceil \frac{n+1}{2} \rceil$ estimates and proposes a new estimate}
 if $p = c_p$ then

```

        wait until [for  $\lceil \frac{n+1}{2} \rceil$  processes  $q$ : received  $(q, r_p, estimate_q, ts_q)$  from  $q$ ]
         $msgs_p[r_p] \leftarrow \{(q, r_p, estimate_q, ts_q) \mid p \text{ received } (q, r_p, estimate_q, ts_q) \text{ from } q\}$ 
         $t \leftarrow \text{largest } ts_q \text{ such that } (q, r_p, estimate_q, ts_q) \in msgs_p[r_p]$ 
         $estimate_p \leftarrow \text{select one } estimate_q \text{ such that } (q, r_p, estimate_q, t) \in msgs_p[r_p]$ 
        send  $(p, r_p, estimate_p)$  to all
    
```

Phase 3: {All processes wait for the new estimate proposed by the current coordinator}

```

    wait until [received  $(c_p, r_p, estimate_{c_p})$  from  $c_p$  or  $c_p \in \mathcal{D}_p$ ] {Query the failure detector}
    if [received  $(c_p, r_p, estimate_{c_p})$  from  $c_p$ ] then           { $p$  received  $estimate_{c_p}$  from  $c_p$ }
         $estimate_p \leftarrow estimate_{c_p}$ 
         $ts_p \leftarrow r_p$ 
        send  $(p, r_p, ack)$  to  $c_p$ 
    else send  $(p, r_p, nack)$  to  $c_p$                                 { $p$  suspects that  $c_p$  crashed}

```

Phase 4: {The current coordinator waits for $\lceil \frac{n+1}{2} \rceil$ replies. If they indicate that $\lceil \frac{n+1}{2} \rceil$ processes adopted its estimate, the coordinator R-broadcasts a decide message}

```

    if  $p = c_p$  then
        wait until [for  $\lceil \frac{n+1}{2} \rceil$  processes  $q$ : received  $(q, r_p, ack)$  or  $(q, r_p, nack)$ ]
        if [for  $\lceil \frac{n+1}{2} \rceil$  processes  $q$ : received  $(q, r_p, ack)$ ] then
            R-broadcast( $p, r_p, estimate_p, decide$ )

```

{If p R-delivers a decide message, p decides accordingly}

```

when R-deliver( $q, r_q, estimate_q, decide$ )
    if  $state_p = \text{undecided}$  then
        decide( $estimate_q$ )
         $state_p \leftarrow \text{decided}$ 

```

Phase 1

Phase 2

Count & Commit.

Termination? Eventually correct p , who no one suspects will become leader.

① S (or w) → Consensus with f+1 nodes

Why gap: ∃ correct process p that is never suspected

↳ Have correct process relay

Challenge: Don't know identity of correct p ahead of time

↳ Agreeing on correct p equivalent to consensus

So instead rely on multiple rounds of communication

Every process p executes the following:

```

procedure propose( $v_p$ )
   $V_p \leftarrow (\perp, \perp, \dots, \perp)$                                 { $p$ 's estimate of the proposed values}
   $V_p[p] \leftarrow v_p$ 
   $\Delta_p \leftarrow V_p$ 

Phase 1: {asynchronous rounds  $r_p$ ,  $1 \leq r_p \leq n - 1$ }
  for  $r_p \leftarrow 1$  to  $n - 1$ 
    send  $(r_p, \Delta_p, p)$  to all
    wait until  $[V_q : \text{received } (r_p, \Delta_q, q) \text{ or } q \in \mathcal{D}_p]$       {query the failure detector}
     $msgs_p[r_p] \leftarrow \{(r_p, \Delta_q, q) \mid \text{received } (r_p, \Delta_q, q)\}$ 
     $\Delta_p \leftarrow (\perp, \perp, \dots, \perp)$ 
    for  $k \leftarrow 1$  to  $n$ 
      if  $V_p[k] = \perp$  and  $\exists(r_p, \Delta_q, q) \in msgs_p[r_p]$  with  $\Delta_q[k] \neq \perp$  then
         $V_p[k] \leftarrow \Delta_q[k]$ 
         $\Delta_p[k] \leftarrow \Delta_q[k]$ 

Phase 2: send  $V_p$  to all
  wait until  $[V_q : \text{received } V_q \text{ or } q \in \mathcal{D}_p]$                       {query the failure detector}
   $lastmsg_p \leftarrow \{V_q \mid \text{received } V_q\}$ 
  for  $k \leftarrow 1$  to  $n$ 
    if  $\exists V_q \in lastmsg_p$  with  $V_q[k] = \perp$  then  $V_p[k] \leftarrow \perp$ 

Phase 3: decide( first non- $\perp$  component of  $V_p$ )

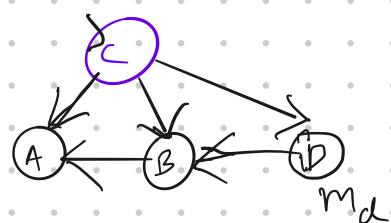
```

Gather & Relay Proposals
Send vector of all proposals

FIG. 5. Solving Consensus using any $\mathcal{D} \in \mathcal{S}$.

Failure Detectors in practice

- Desirable but
 - Accuracy is hard



↳ Weak: How to ensure that some process p is never suspected?

↳ How to choose p ?

- Strong: No correct process is even suspected
- Delay vs failure
- Prior work (Falcon, others)
 - Kill suspected nodes
 - Likely to violate any assumptions about # of process failures.

Pigeon

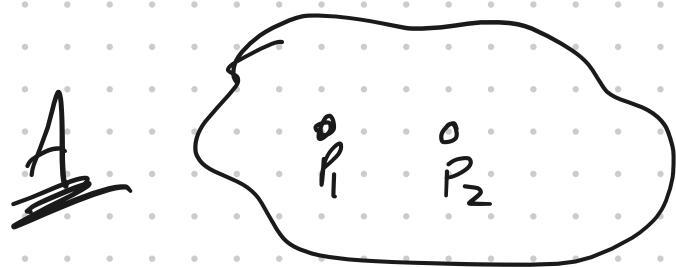
→ Provide information about failure certainty

(Warnings vs facts)

→ Use additional sources of information to] Core Benefit!
improve predictions

Strong Accuracy

Weak Complete



Output

$$\hookrightarrow f_d(x) = \{P_1, P_2\}$$

- Can P_1 be correct?
- Can P_1 have failed?

SA/WC $\xrightarrow{\text{Union}}$ St/SC

WA

WC [Sp. Process
C]

A



A

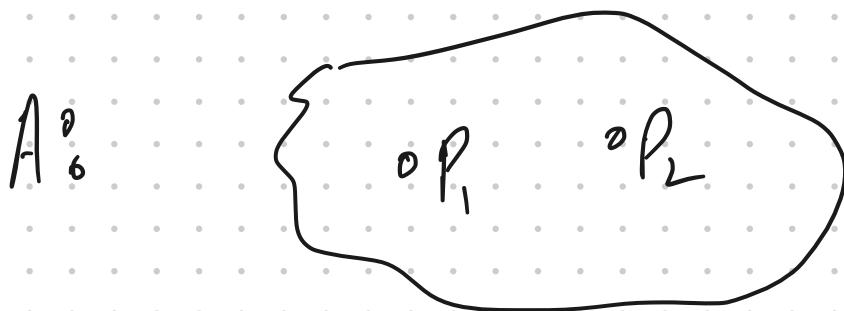
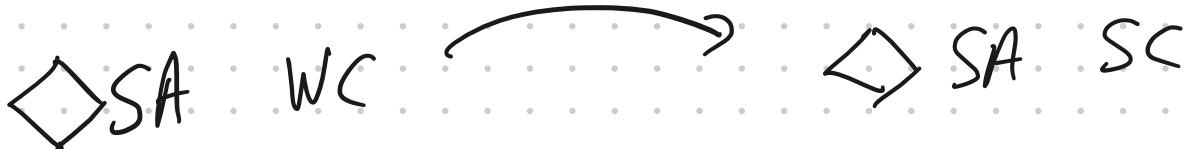
Can P_1 be correct?

Yes

Is P_2 failed? Yes

② Can P_1 be faulty?

③ Can $P_1 = C$?



Q1. Can P_1 be correct?

Q2. Can P_1 be faulty?

A wavy line with a circle containing a question mark, followed by the text "No correct process is suspected".

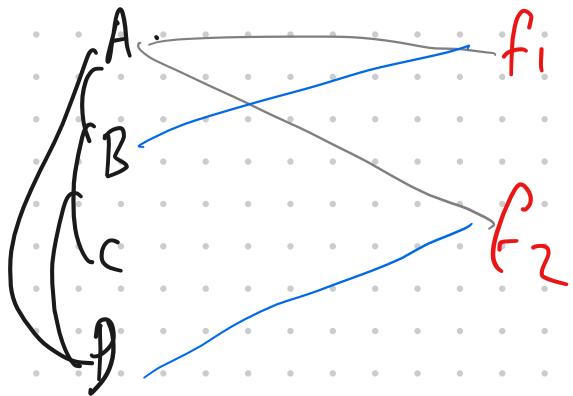
*No correct process
is suspected*

Do not know

Weak Complete

Correct

Failed



Eventually

Strong Complete

A

f_1

B

f_2

C

f_3

D

f_4

B

C

D

f_1

f_2

f_3

f_4