# Circumventing Impossibility

Failure Detectors

# Circumventing Impossibility

- Use timing assumptions
  - $\Delta$, $\Phi$
- Timing assumptions can be difficult to work with
- What about time-free algorithms?

# Time-Free Algorithms

- Describe an algorithm using a *failure detector* abstraction [Chandra & Toueg 96]
- By abstracting away time we
  - get simpler algorithms
  - could specify (minimum) conditions under which certain problems are solvable in a simpler way
  - characterize systems based on their ability to implement certain types of failure detectors

# Environment Model

- Asynchronous message-passing system
  - fully connected
- Crash failures
- n process $P_1,\ldots,P_n$
- Reliable links between each pair of correct processes
  - Can implement Reliable Broadcast

# The Failure Detector Abstraction

- Each process has local failure detector oracle
  - typically outputs list of processes suspected to have crashed at any given time
- In each execution step, a process
  - receives a message (if there is one ready)
  - queries its failure detector oracle
  - makes a transition to a new state
  - may send messages to other processes

# Specifying Failure Detectors

- *Accuracy*: which processes are not suspected and when
- *Completeness*: which processes are suspected and when

- Why do we need both?
  - A trivial failure detector will satisfy any accuracy (completeness) property by never (always) suspecting any (all) 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 process is suspected before it crashes
- *Weak Accuracy*: Some correct process is never suspected
- *Eventual Strong Accuracy*: Eventually correct processes are not suspected by any correct process
- *Eventual Weak Accuracy*: Eventually some correct process is never suspected by any correct process

# Failure Detector Classes

| Completeness | Accuracy | | | |
|---|---|---|---|---|
| | Strong | Weak | Eventual Strong | Eventual Weak |
| Strong | Perfect<br>P | Strong<br>S | Eventually Perfect<br>◊P | Eventually Strong<br>◊S |
| Weak | Q | Weak<br>W | ◊Q | Eventually Weak<br>◊W |

# Reducibility of failure detectors

- If a f. d. D' can be implemented using a f. d. D, then
  - D ≥ D', or D' is *reducible* to D, or D' is *weaker* than D
- If D ≥ D' and D'≥ D ➔ D and D' are equiv.

# Reducibility of classes of failure detectors

- $C$ and $C'$ are classes of f. d. and $\forall D \in C, \exists D' \in C'$: $D \geq D' \Rightarrow C \geq C'$
  - $C \geq C'$, or $C'$ is weaker than $C$
- If $C \geq C'$ and $C' \geq C \Rightarrow C$ and $C'$ are equiv.
- If $C \geq C'$, then any problem solvable using a $D' \in C'$ is also solvable using some $D \in C$

- $P \geq Q,\ S \geq W,\ \Diamond P \geq \Diamond Q,\ \Diamond S \geq \Diamond W$
- $P > S,\ \Diamond P > \Diamond S,\ P > \Diamond P,\ S > \Diamond S,\ P > \Diamond S$
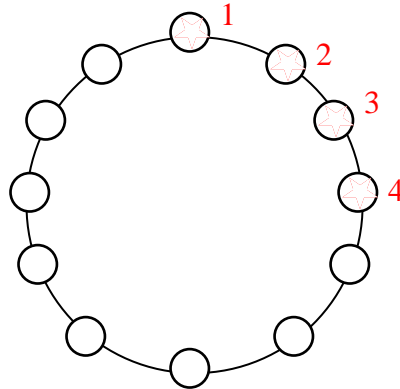
# Weak Completeness ➔ Strong Completeness

- Let D be a failure detector satisfying Weak Completeness
- Construct a failure detector D' satisfying Strong Completeness
- Conclusion: $D \geq D'$
- Theorem: $Q \geq P,\ W \geq S,\ \Diamond Q \geq \Diamond P,\ \Diamond W \geq S$
- Corollary: $Q = P,\ W = S,\ \Diamond Q = \Diamond P,\ \Diamond W = \Diamond S$

## Solving Consensus using $\Diamond$ S : Rotating Coordinator Algorithms

Work for up to $f < n/2$ crashes

• Processes are numbered 1, 2, …, n

• They execute asynchronous *rounds*

• In round r , the *coordinator* is
  process (r mod n) + 1

• In round r , the coordinator:
  - tries to impose its estimate as the consensus value
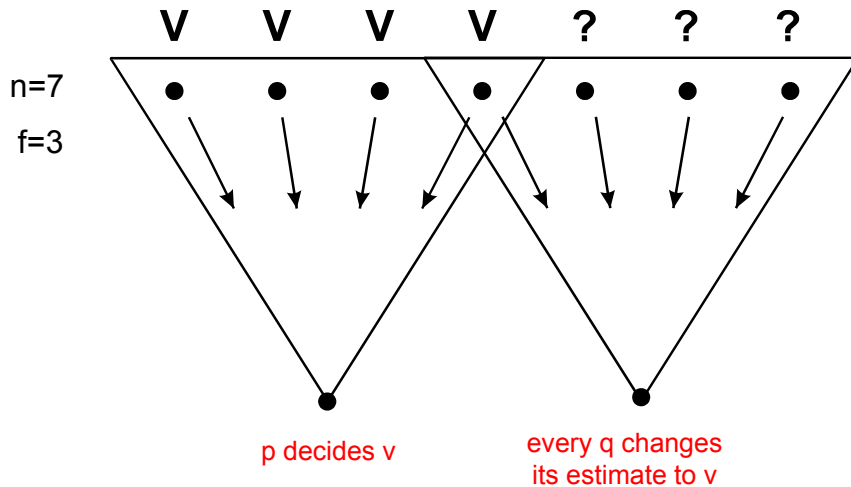  - succeeds if does not crash and it is not suspected by $\Diamond$ S

---

## Consensus using $\Diamond$S : Phase k; p = k mod n

• q sends (q, k, estimate, $r_{estimate}$) to p
  - p awaits a majority of round k estimates and chooses the estimate EST with the highest $r_{estimate}$
• p broadcasts (p,EST,k)
• q waits to hear (p,EST,k), or FD suspects p
  - if q hears from p, then q responds with (ack, k) and estimate := EST; $r_{estimate}$ := k
  - else q sends (nack, k)
• p awaits a majority of round k ack/nack's
  - If acks, then p R-broadcasts (decide v)

# Why does it work?

Agreement:



n=7

f=3

V    V    V    V    ?    ?    ?

p decides v

every q changes
its estimate to v

---

# Agreement

- Lemma 1: Let r be the smallest round such that $c = r \bmod n$ decides. Let $est_c$ be the estimate sent by c prior to deciding. Then

  For all rounds $r' \geq r$: If $c' = r' \bmod n$ sends $est_{c'}$, then $est_{c'} = est_c$

Proof: By induction on r'. Trivial for r'=r. Assume for r≤r'<k. Prove for r'=k. Since a majority responds to c with ack, one of the responding processes p sent its estimate to $c_k$. Hence $ts_p \geq r$. Since $c_k$ selects maximum, it will select the estimate of q with timestamp t, r ≤t<k. By the induction hypothesis $est_q = est_c$.

# Consensus using ◊S

- Agreement is immediate from Lemma 1
- Termination:
  - No process blocks indefinitely in any one of the wait statements. Why?
  - Eventually, all processes do not suspect a single correct process. Once it becomes a coordinator, everybody decides
- Validity: obvious

# Consensus using ◊S

- Theorem: Consensus among n>2t is solvable using ◊S
- Corollary: Consensus among n>2t is solvable using ◊W
  - Why?
  - because ◊W ≥ ◊S
- The bound on the number of processes is tight (partitioning argument)

# Implementation with partial synchrony

- Assume $\Delta$ is unknown
- Each process p holds a local estimate $\Delta_p$
- Broadcast "I am alive" to all every time unit
- If p does not receive "I am alive" from q during the last $\Delta_p$ time units → suspect q
- If p receives a message from q and q is suspected→ stop suspecting q, $\Delta_p := \Delta_p + 1$

To which class belongs the failure detector implemented above?