# Coordination and Agreement

# Outline

- Introduction

- Distributed Mutual Exclusion

- Election Algorithms

- Group Communication

- Consensus and Related Problems

# Introduction
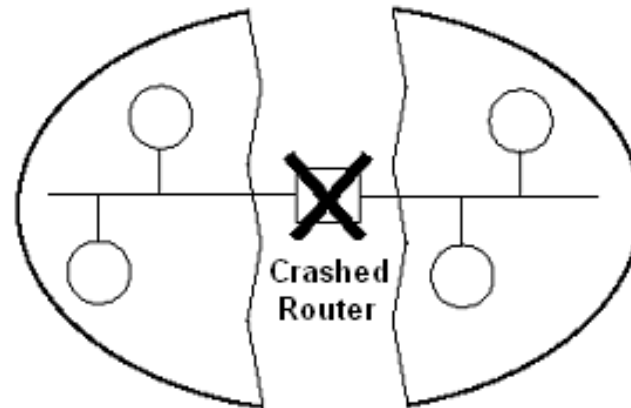
- Collection of algorithms that share an aim

  for a set of processes to coordinate their actions or to agree on one or more values.

- For example, Spaceship:

  - **Agreement:** it is essential that the computers controlling it agree on such conditions as whether the spaceship's mission is proceeding or has been aborted

  - **Coordination:** the computers must coordinate their actions correctly with respect to shared resources (the spaceship's sensors and actuators)

# Main Assumptions

- Each pair of processes is connected by reliable channels

- Processes independent from each other

- Network: don't disconnect


Crashed Router

- Processes fail only by crashing

- Local failure detector

# Failure Detector

- Is a service that processes queries about whether a particular process has crashed.

- It is often implemented by a local object known as a *Local Failure Detector*.

- Failure detectors are not necessarily accurate.

- For example:
  - a process that timed-out after 255 seconds might have succeeded if allowed to proceed for 256 seconds.

- Two types of failure detector:
  - Unreliable  failure detector
  - Reliable  failure detector

# Unreliable Failure Detector

- Produce one of two values when given the identity of a process: *Unsuspected* or *Suspected*.

  - **Unsuspected:** detector has recently received evidence suggesting that the process has not failed.

  - **Suspected:** failure detector has some indication that the process may have failed.

- Implement:

  - each process sends *alive* message to everyone else

  - not receiving *alive* message after timeout, report *Suspected*

  - if it subsequently receives, reports OK *(Unsuspected)*
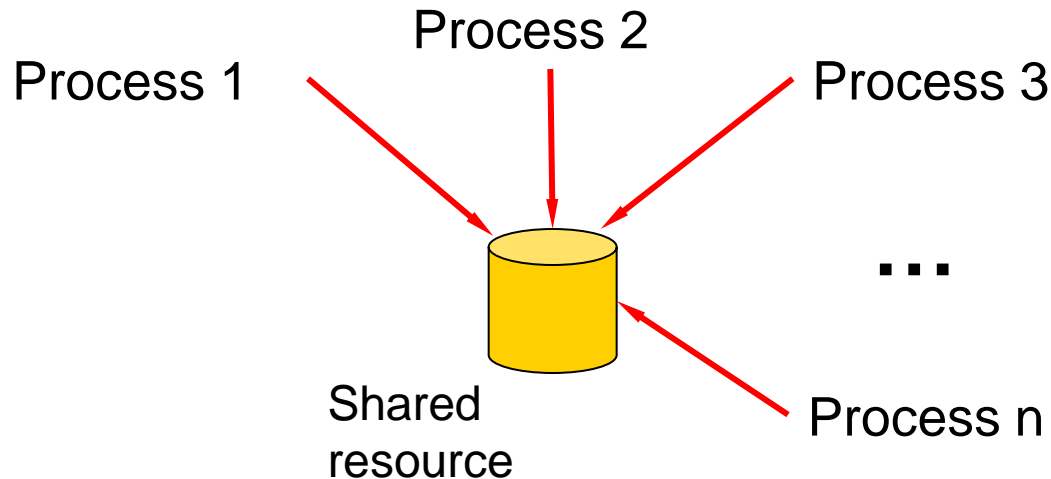
- Most practical systems

# Reliable Failure Detector

- Is always accurate in detecting a process's failure.
- It answers processes' queries with either a response of *Unsuspected* or *Failed*.
  - **Unsuspected:** as before, can only be a hint that the process has not failed.
  - **Failed:** detector has determined that the process has crashed.
- Implement needs *synchronous* system
- Few practical systems

# Outline

- Introduction

- **Distributed Mutual Exclusion**

- Election Algorithms

- Group Communication

- Consensus and Related Problems

# Distributed Mutual Exclusion (1)

Process 2

Process 1

Process 3

...

Shared
resource

Process n

- Mutual exclusion very important
  - Prevent interference
  - Ensure consistency when accessing the resources

# Distributed Mutual Exclusion (2)

- Mutual exclusion useful when the server managing the resources don't use locks

- Critical section

*Enter()*          *enter critical section – blocking*

   •
   •          *Access shared resources  in critical*
   •          *section*

*Exit()*          *Leave critical section*

# Distributed Mutual Exclusion (3)

- Distributed mutual exclusion: no shared variables, only message passing

- Properties:

    - **Safety:** At most one process may execute in the critical section at a time

    - **Liveness:** Requests to enter and exit the critical section eventually succeed

        No deadlock and no starvation

    - **Ordering:** If one request to enter the CS happened-before another, then entry to the CS is granted in that order
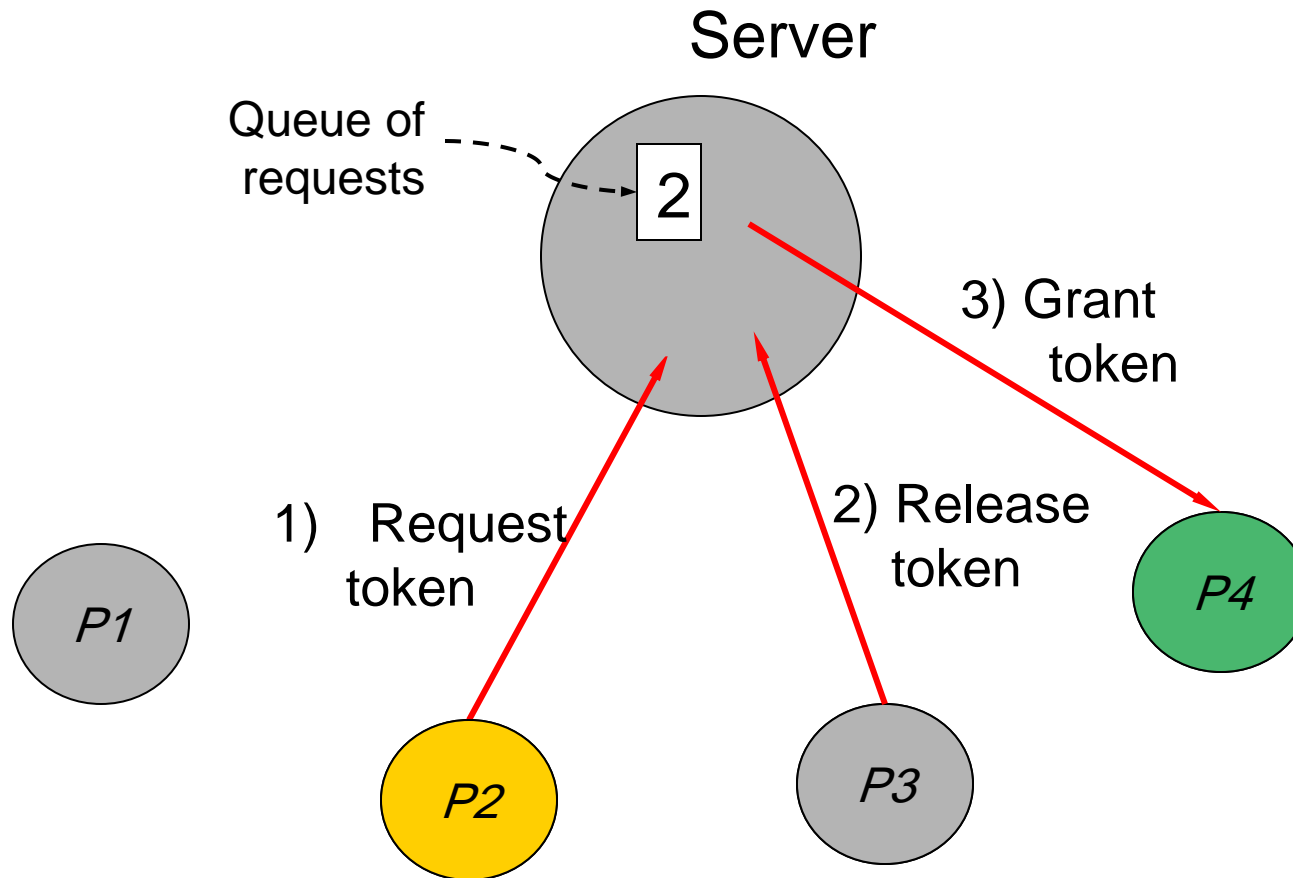
# Mutual Exclusion Algorithms

- Basic Hypothesis:
  - System: asynchronous
  - Critical section: only one
  - Processes: don't fail
  - Message transmission: reliable
- Central Server Algorithm
- Ring-Based Algorithm
- Mutual Exclusion using Multicast and Logical Clocks
- Maekawa's Voting Algorithm
- Mutual Exclusion Algorithms Comparison

# Evaluation of the performance alg.

- Bandwidth
  - The number of message sent in each entry and exit operation
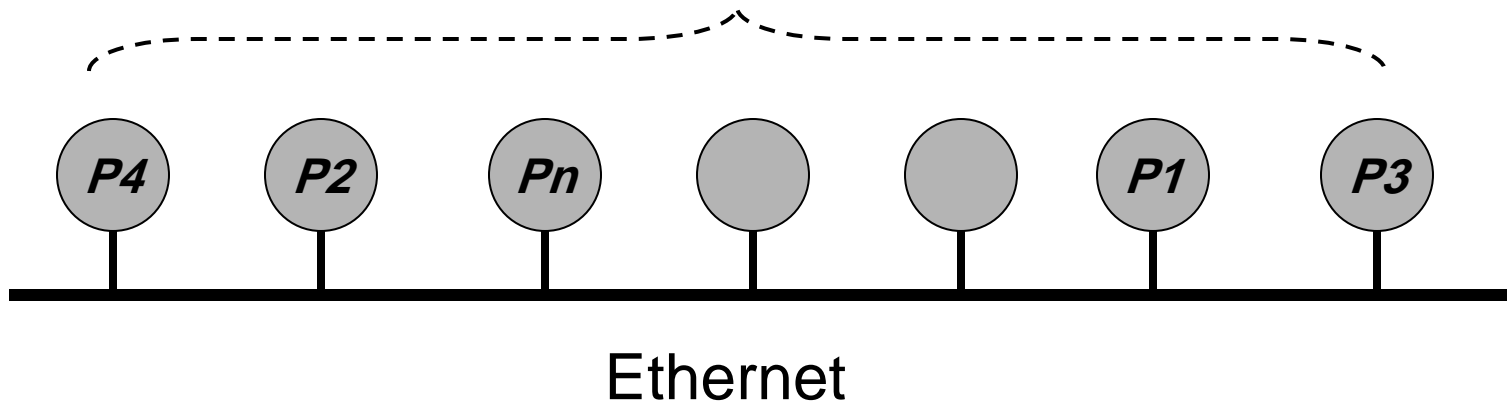
- Client Delay

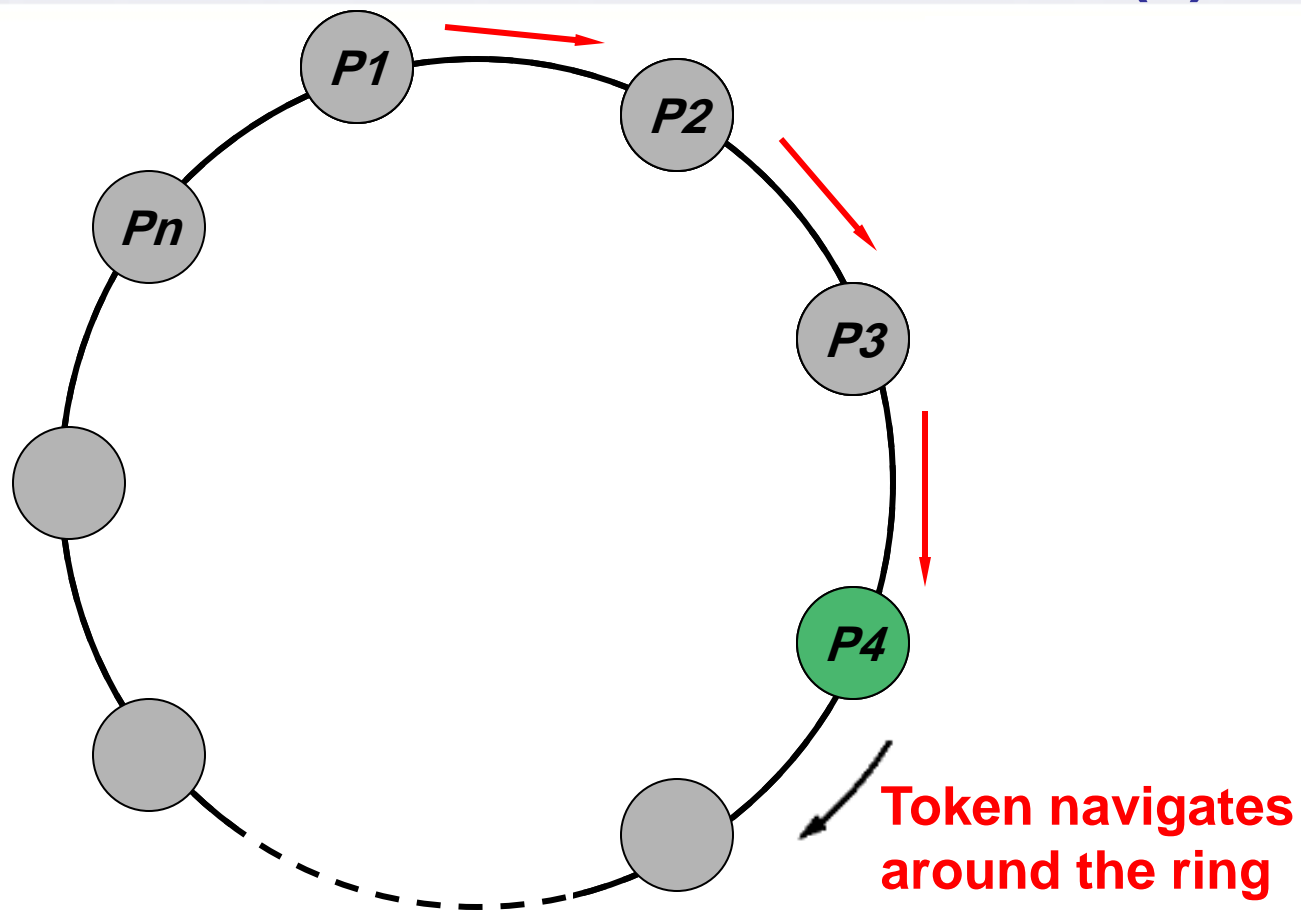- Throughput

# Central Server Algorithm

# Ring-Based Algorithm (1)

A group of unordered
processes in a network



P4  P2  Pn  ( )  ( )  P1  P3

Ethernet

# Ring-Based Algorithm (2)



Token navigates around the ring

# Mutual Exclusion using Multicast and Logical Clocks



**Waiting queue**

P3

19

19

OK

23

OK

P1

OK

23

OK

19

P2

23

P1 and P2 request entering the critical section simultaneously

- Main steps of the algorithm:

**Initialization**

**State := RELEASED;**

**Process $p_i$ request entering the critical section**

**State := WANTED;**

**T := request's timestamp;**

**Multicast request <T, $p_i$> to all processes;**

**Wait until (Number of replies received = ($N - 1$));**

**State := HELD;**

# **Mutual Exclusion using Multicast and Logical Clocks** (3)

- Main steps of the algorithm (cont'd):

**On receipt of a request <$T_i$, $p_i$> at $p_j$ (i $\neq$ j)**

**If (state = HELD) OR**

**(state = WANTED AND (T, $p_j$) < ($T_i$, $p_i$))**

**Then queue request from $p_i$ without replying;**

**Else reply immediately to $p_i$;**

**To quit the critical section**

**state := RELEASED;**

**Reply to any queued requests;**

# Maekawa's Voting Algorithm (1)

- Candidate process: must collect sufficient votes to enter to the critical section

- Each process $p_i$ maintain a *voting set* $V_i$ (i=1, ..., N), where $V_i \subseteq \{p_1, ..., p_N\}$

- Sets $V_i$: chosen such that $\forall$ i,j

  - $p_i \in V_i$

  - $V_i \cap V_j \neq \varnothing$     **(at least one common member of any two voting sets)**

  - $|V_i| = k$     **(fairness)**

  - Each process $p_j$ is contained in M of the voting sets $V_i$

# Maekawa's Voting Algorithm (2)

- Main steps of the algorithm:

**Initialization**

**state := RELEASED;**

**voted := FALSE;**

**For $p_i$ to enter the critical section**

**state := WANTED;**

**Multicast request to all processes in $V_i$ ;**

**Wait until (number of replies received = $K$);**

**state := HELD;** <span style="color:red">**$p_i$ enter the critical section only after collecting K votes**</span>

# Maekawa's Voting Algorithm (3)

- Main steps of the algorithm (cont'd):

**On receipt of a request from $p_i$ at $p_j$**

**If (state = HELD OR voted = TRUE)**

**Then** queue request from $p_i$ without replying;

**Else** Reply immediately to $p_i$;
voted := TRUE;

**For $p_i$ to exit the critical section**

**state := RELEASED;**

**Multicast release to all processes $V_i$ ;**

- Main steps of the algorithm (cont'd):

**On a receipt of a release from $p_i$ at $p_j$**

**If (queue of requests is non-empty)**

    **Then**     **remove head of queue, e.g., $p_k$;**

            **send reply to $p_k$;**

            **voted := TRUE;**

    **Else voted := FALSE;**

# M. E. Algorithms Comparison

| Algorithm | Number of messages | | Problems |
|---|---|---|---|
| | Enter()/Exit | Before Enter() | |
| Centralized | 3 | 2 | Crash of server |
| Virtual ring | 1 to N | 0 to N-1 | Crash of a process Token lost Ordering non-satisfied |
| Logical clocks | 3(N-1) | 2(N-1) | Crash of a process |
| Maekawa's Alg. | $3\sqrt{N}$ | $2\sqrt{N}$ | Crash of a process who **votes** |

# Outline

- Introduction

- Distributed Mutual Exclusion

- **Election Algorithms**

- Group Communication

- Consensus and Related Problems

# **Election Algorithms** (1)

- **Objective:** Elect one process $p_i$ from a group of processes $p_1 \ldots p_N$

- At any point in time, a process $p_i$ is either a participant or a non-participant
  
  **Even if multiple elections have been started simultaneously**

- Each process $p_i$ maintains the identity of the elected in the variable *Elected*$_i$ (NIL '⊥' if it isn't defined yet)

- **Properties to satisfy:** ∀ $p_i$

  - Safety**:** *Elected*$_i$ = NIL or *Elected* = P

    **A non-crashed process with the largest identifier**

  - Liveness**:** $p_i$ participates and sets *Elected*$_i$ ≠ NIL, or crashes

# Election Algorithms (2)

- Ring-Based Election Algorithm

- Bully Algorithm

- Election Algorithms Comparison

**Process 5 starts the election**

# Ring-Based Election Algorithm (2)

**Initialization**

> $Participant_i := FALSE;$
> $Elected_i := NIL$

**$P_i$ starts an election**

> $Participant_i := TRUE;$
>
> Send the message *<election, $p_i$>* to its neighbor

**Receipt of a message *<elected, $p_j$> at $p_i$***

> **If** $p_i \neq p_j$
> **Then** $Participant_i := FALSE;$
> $Elected_i := p_j;$
> Send the message *<elected, $p_j$>* to its neighbor

**Receipt of the election's message $<election, p_i>$ at $p_j$**

**<u>If</u>**  $p_i > p_j$

**<u>Then</u>**     Send the message $<election, p_i>$ to its neighbor
              $Participant_j := TRUE$;

**<u>Else If</u>** $p_i < p_j$ AND $Participant_j = FALSE$

**<u>Then</u>**     Send the message $<election, p_i>$ to its neighbor
              $Participant_j := TRUE$;

**<u>Else If</u>**  $p_i = p_j$

**<u>Then</u>**     $Elected_j := p_j$;
              $Participant_j := FALSE$;
              Send the message $<elected, p_j>$ to its neighbor

# Bully Algorithm (1)

- **Characteristic:** Allows processes to crash during an election

- **Hypothesis:**
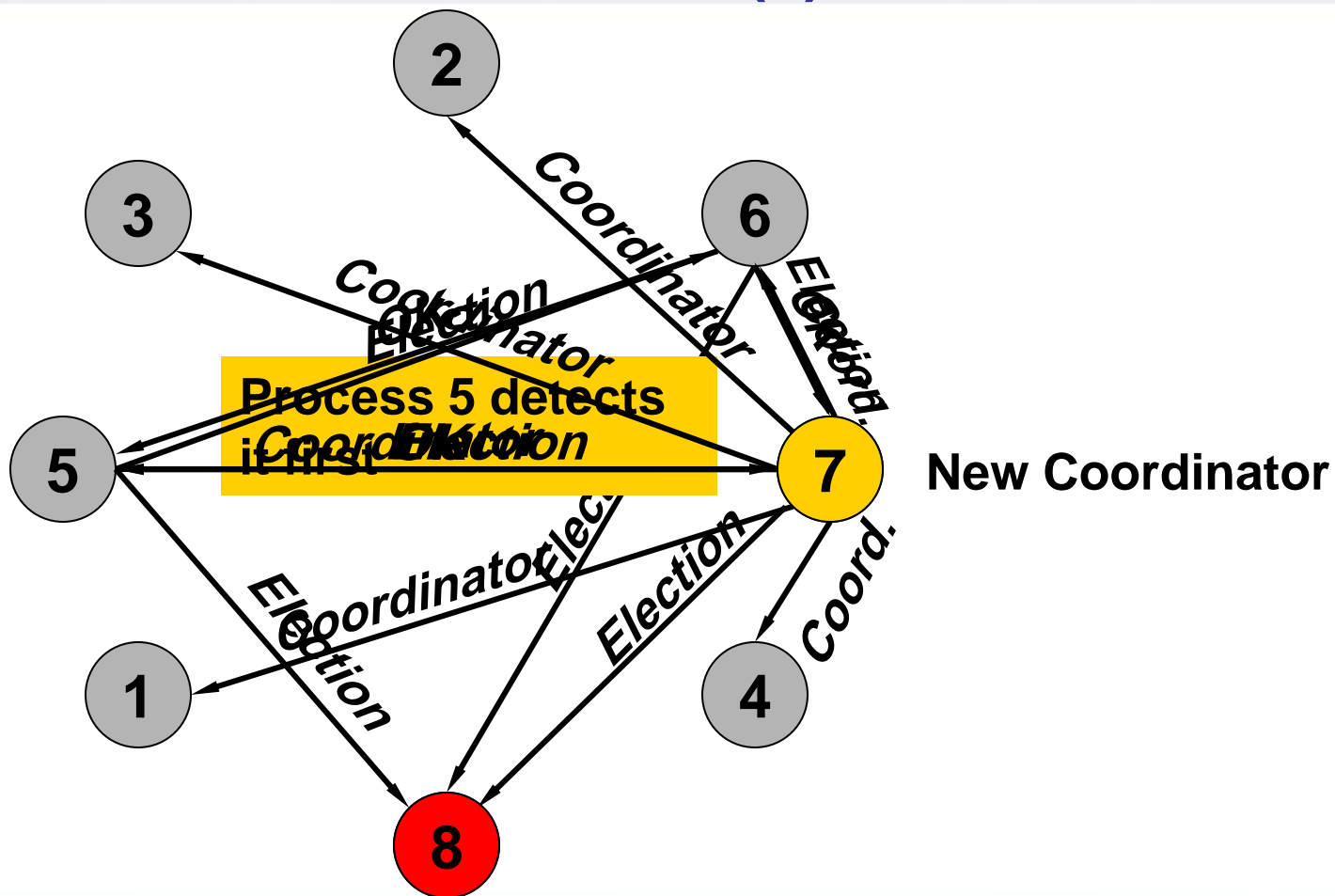
  - Reliable transmission

  - Synchronous system

$Delay_{Trans.}$

$Delay_{Trans.}$

$Delay_{Trait.}$

$$T = 2\ Delay_{Trans.}\ +\ Delay_{Trait.}$$

# Bully Algorithm (2)

- **Hypothesis (cont'd):**
  - Each process knows which processes have higher identifiers, and it can communicate with all such processes

- **Three types of messages:**
  - *Election*: starts an election
  - *OK*: sent in response to an election message
  - *Coordinator*: announces the new coordinator

- Election started by a process when it notices, through timeouts, that the coordinator has failed

# Bully Algorithm (3)



**Process 5 detects it first**

**New Coordinator**

33

# Bully Algorithm (4)

**Initialization**

Elected$_i$ := NIL

**p$_i$ starts the election**

Send the message (*Election, p$_i$*) to p$_j$ , i.e., p$_j$ > p$_i$

Waits until message (*OK, p$_j$*) from p$_j$ are received;

<u>If</u> no message (*OK, p$_j$*) arrives during T

<u>Then</u>    Elected$_i$ := p$_i$;
             Send the message (*Coordinator, p$_i$*) to p$_j$ , i.e., p$_j$ < p$_i$

<u>Else</u>  waits until receipt of the message (*coordinator*)
(if it doesn't arrive during  another timeout T', it begins another election)

# Bully Algorithm (5)

**Receipt of the message (*Coordinator, $p_j$*) at $p_i$**

**Elected$_i$ := p$_j$;**

**Receipt of the message (*Election, $p_j$* ) at $p_i$**

**Send the message *(OK, $p_i$)* to p$_j$**

**Start the election unless it has begun one already**

- When a process is started to replace a crashed process: it begins an election

# Election Algorithms Comparison

| Election algorithm | Number of messages | Problems |
|---|---|---|
| Virtual ring | 2N to 3N-1 | Don't tolerate faults |
| Bully | N-2 to $O(N^2)$ | System must be synchronous |