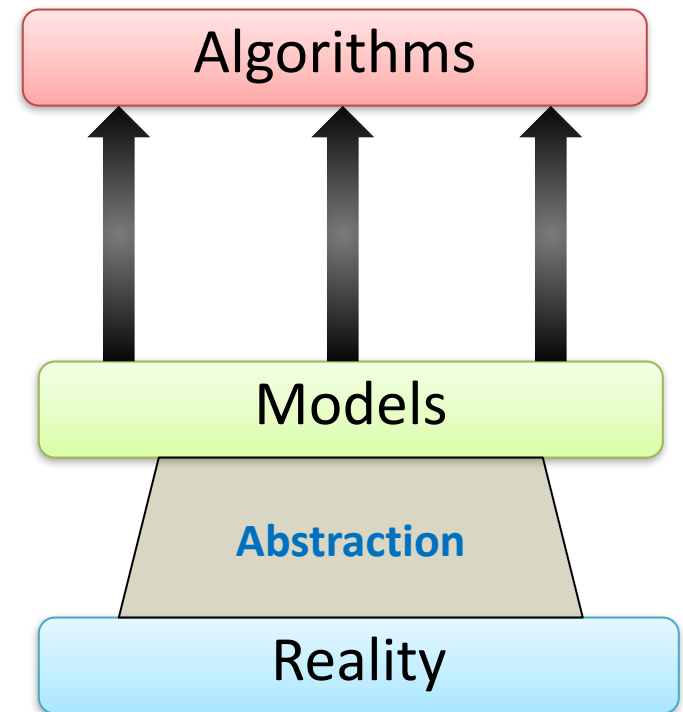Pixabay license

# DISTRIBUTED SYSTEM MODELS

# Distributed system model
## Of theoretical relevance for designing algorithms

- Model captures all the **assumptions** made about the system

- Including network, clock, processor, etc.

- Algorithms always **assume** a certain model

- Some algorithms are only possible in **stronger models** (more restrictions)

- Model is of **theoretical relevance** whether its assumptions hold in practice is a different question

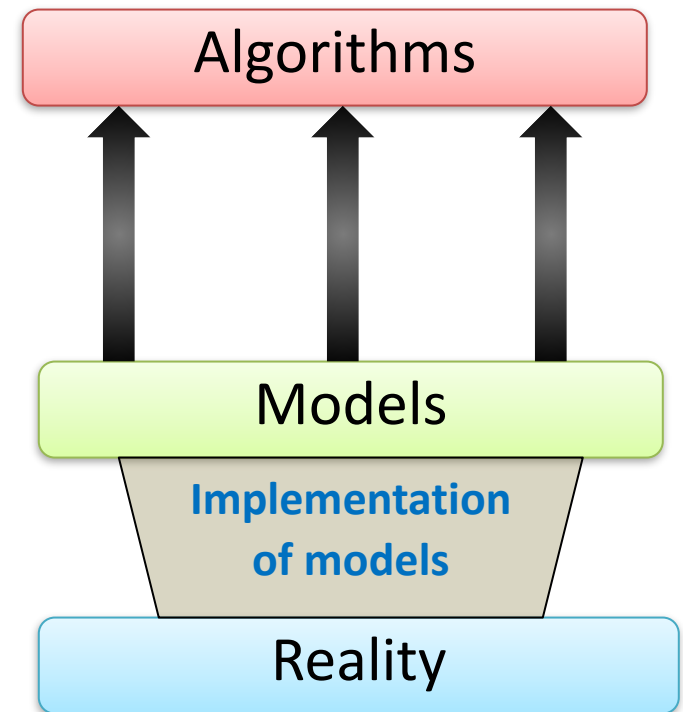*"All models are wrong, but some are useful"*

George Box, ~1976

Algorithms

Models

Abstraction

Reality

# Distributed system model
## Of theoretical relevance for designing algorithms

- Model captures all the **assumptions** made about the system

- Including network, clock, processor, etc.

- Algorithms always **assume** a certain model

- Some algorithms are only possible in **stronger models** (more restrictions)

- Model is of **theoretical relevance** whether its assumptions hold in practice is a different question

*"All models are wrong, but some are useful"*

George Box, ~1976

**Algorithms**

**Models**

**Implementation of models**

**Reality**

# Synchronous vs. asynchronous model

| Property | Synchronous system model | Asynchronous system model |
|---|---|---|
| Clocks | Bound on drift | No bound on drift |
| Processor | Bound on execution time | No bound on execution time |
| Channel | Bound on latency | No bound on latency |

# Synchronous distributed system model

- Each process has a local clock whose drift rate from real time is **within a known bound**

- Each step of a computation will complete **within a known bound**

- Each message transmitted over a channel is received **within a known bound**

# Asynchronous distributed system model

- Clock **drift rates are arbitrary**

- Each step of a computation may take an **arbitrary time to complete** (but will eventually complete)

- Messages may need an **arbitrary time to be transmitted** (but will eventually be transmitted)

# Two General's Problem (Agreement)

Armies need to agree on:

Army
Part 1

Army
Part 2

Enemy

**Armies are safe if they don't attack (or win) (Safety)**

**Armies need to coordinate attack to win. (Liveness)**

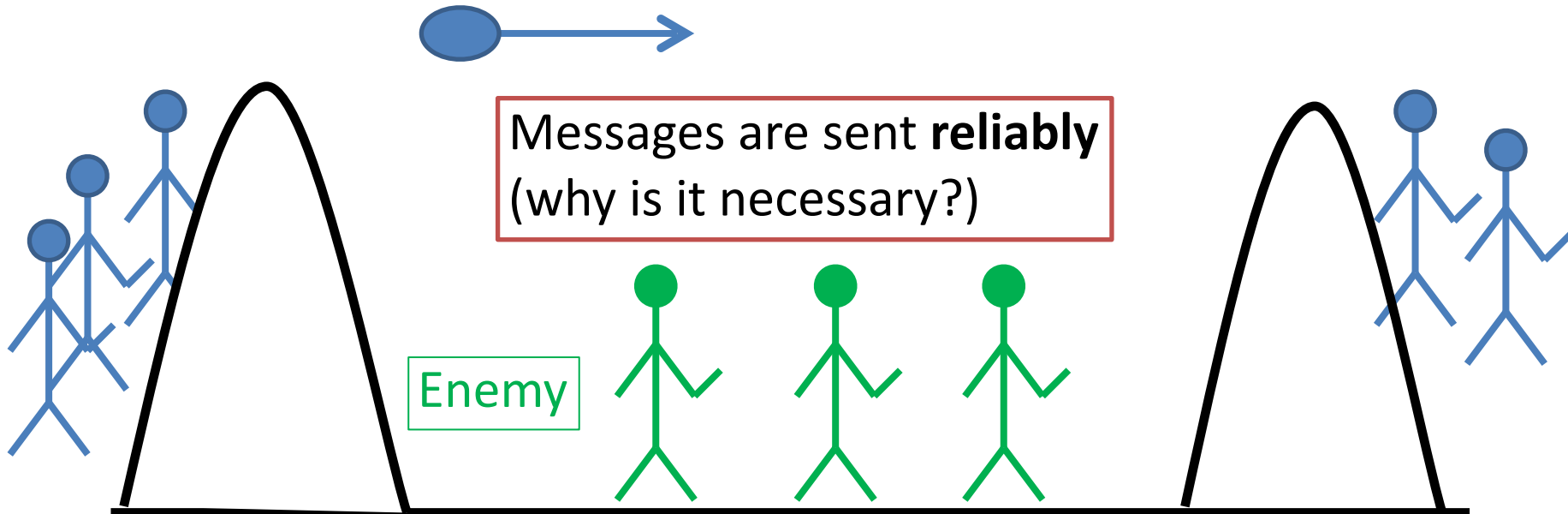# Two General's Problem (Agreement)

Armies need to agree on:

*Who leads the attack?*
*When to attack?*

A thought experiment

Army
Part 1

Army
Part 2



Enemy

Armies are safe if they don't attack (or win) (Safety)

Armies need to coordinate attack to win. (Liveness)

# Two General's Problem (Agreement)

Armies need to agree on:

*Who leads the attack?*
*When to attack?*

A thought experiment

Army
Part 1

Army
Part 2

Messages are sent **reliably**
(why is it necessary?)

Enemy

Armies are safe if they
don't attack (or win)
(Safety)

Armies need to
coordinate attack to
win.
(Liveness)

# Synchronous vs. asynchronous agreement

*Who leads the attack?*
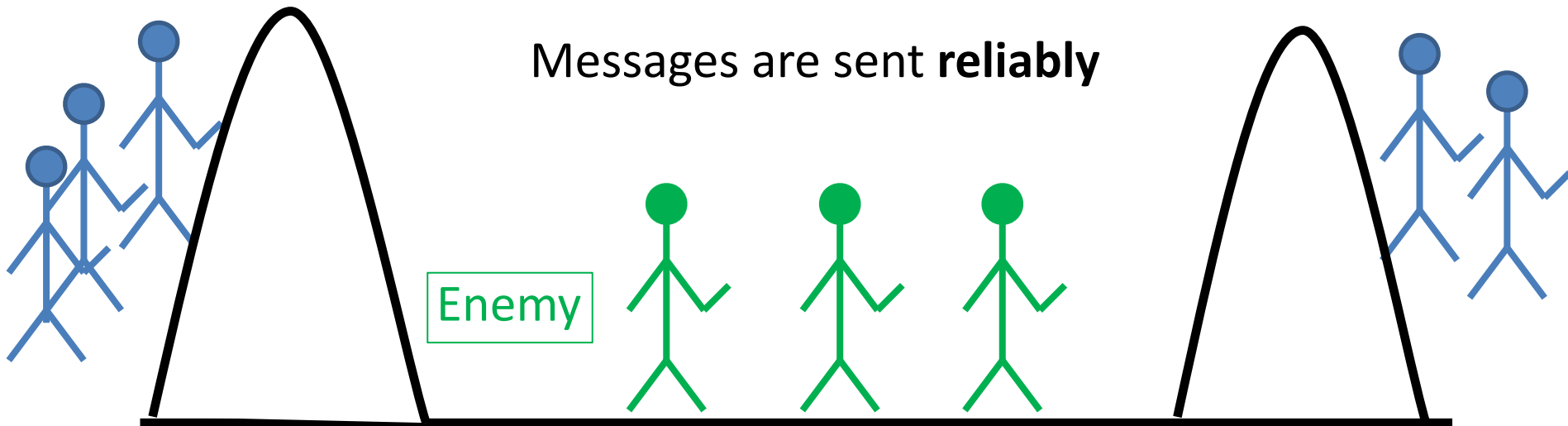
- Largest army
- If tied, Army 1

Army
Part 1

Army
Part 2

**2**

**3**

Messages are sent **reliably**

Enemy

# Asynchronous agreement

*When to attack? No bound on delivery!*

Army
Part 1

Army
Part 2

Messages are sent **reliably**

Enemy

# Asynchronous agreement

*When to attack? No bound on delivery!*

Army
Part 1

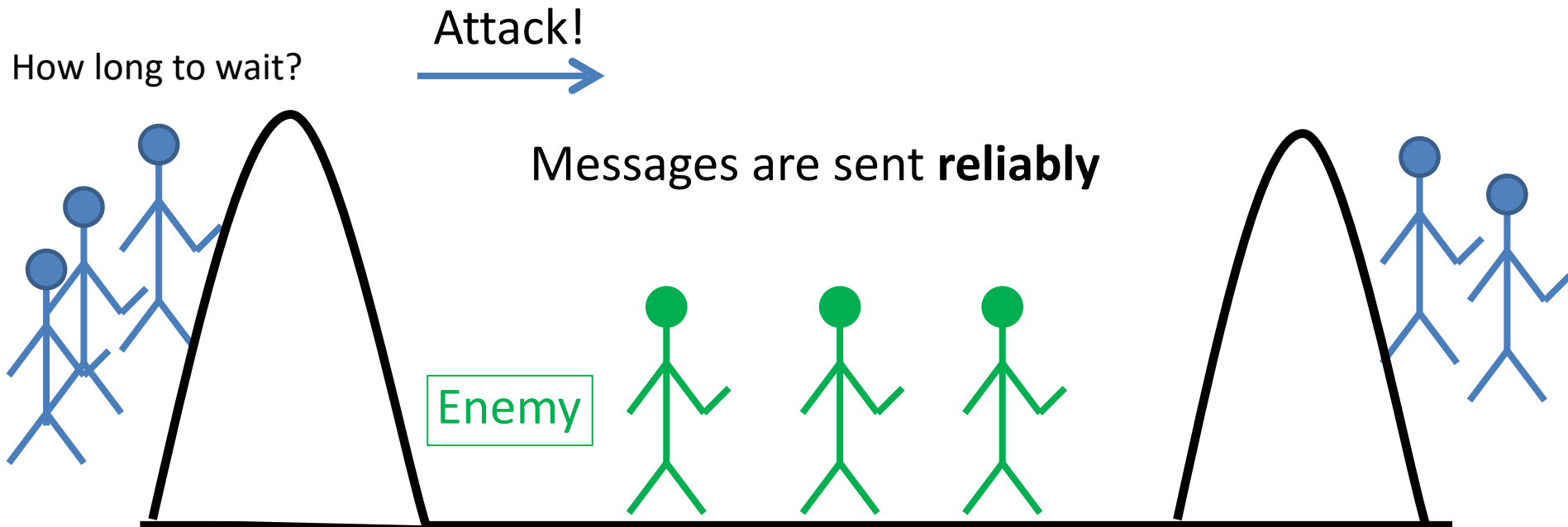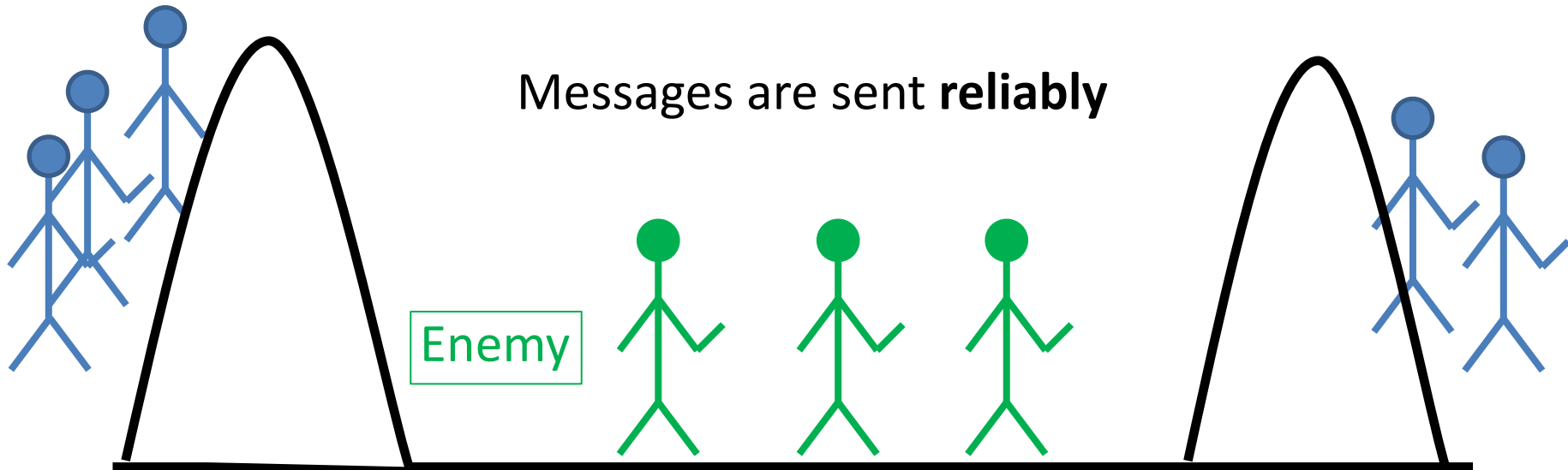Army
Part 2

Attack!

Messages are sent **reliably**

Enemy

# Asynchronous agreement

*When to attack? No bound on delivery!*

Army
Part 1

Army
Part 2

How long to wait?

Attack!

Messages are sent **reliably**

Enemy

# Synchronous agreement

## *When to attack?*

Army
Part 1

Army
Part 2

Messages are sent **reliably**

Enemy

# Synchronous agreement
*When to attack?*

Army
Part 1

Attack!

Army
Part 2

Messages are sent **reliably**

Enemy

# Synchronous agreement

*When to attack?*

Message takes at least ***min*** time and at most ***max*** time to arrive

Army
Part 1

Attack!
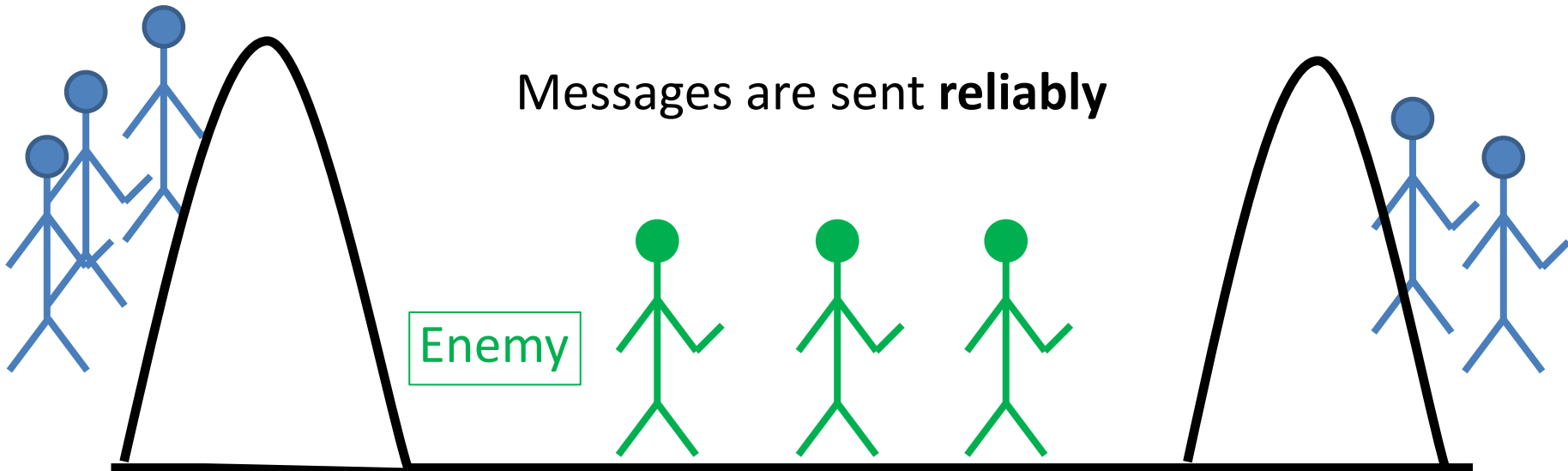
Army
Part 2

Messages are sent **reliably**

Enemy

# Synchronous agreement
## *When to attack?*

Message takes at least **min** time
and at most **max** time to arrive

Attack!

**Strategy**: Waits for *min* time then attacks.
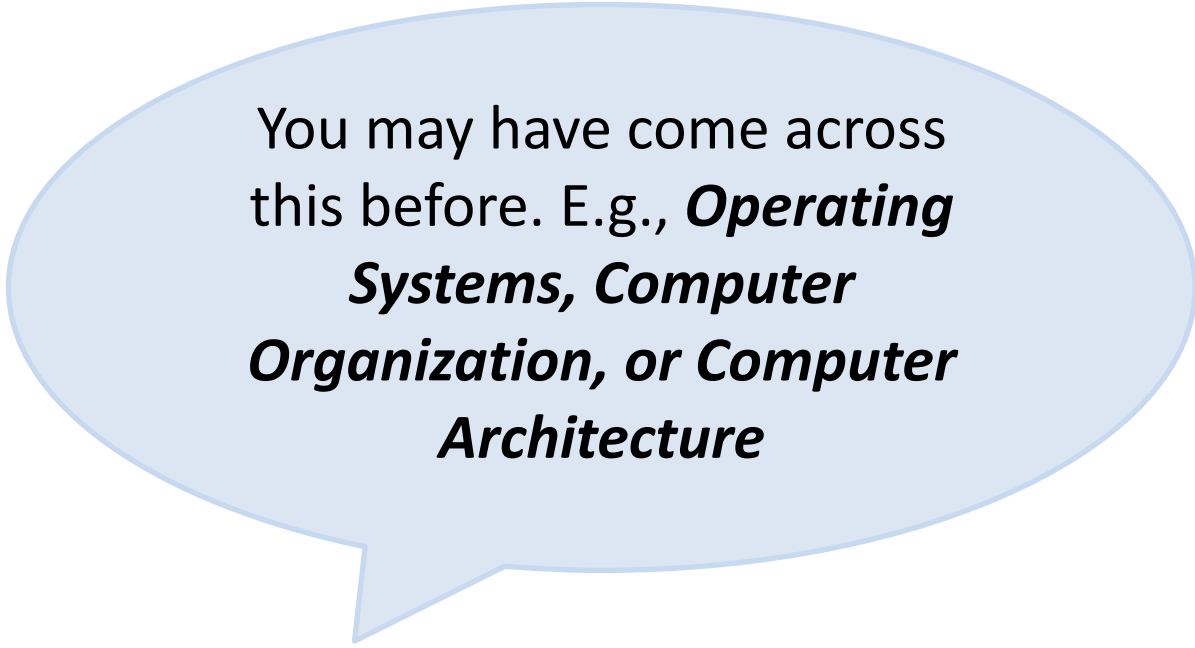
Messages are sent **reliably**

Army
Part 1

Army
Part 2

Enemy

# Synchronous agreement

## *When to attack?*

Message takes at least **min** time
and at most **max** time to arrive

Army
Part 1

Army
Part 2

Attack!

→

**Strategy**: Waits for *min* time then attacks.

Messages are sent **reliably**

Enemy

**Guarantee:** Army 2 attacks
no later than **max − min**
time after Army 1.

# **Summarizing takeaways**

- Some problems cannot be solved in an asynchronous world (e.g., *when* vs. *who* leads attack)

- A solution valid for asynchronous distributed systems is also valid for synchronous ones (synchronous model is a stronger model)

- Internet and many practical distributed applications are closer to asynchronous than synchronous model

- Apply timeouts and timing assumptions to reduce uncertainty and to bring elements of the synchronous model into the picture

# Self-study questions

- Think of a few design problems that cannot be solved in an asynchronous world …

- … and show how they can be solved in a synchronous world.

- What are some other useful assumptions for distributed system design and how practical are they?

- Can you think of other coordination and agreement problems?

You may have come across this before. E.g., *Operating Systems, Computer Organization, or Computer Architecture*
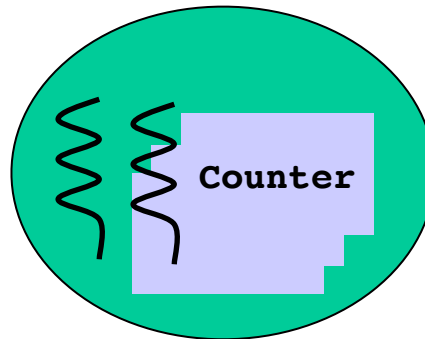
# RECAP: MUTUAL EXCLUSION

## IN A NON-DISTRIBUTED SYSTEMS CONTEXT

# Problem: Access to shared variables

- Imagine a **globally shared variable** `counter` accessible to multiple threads (or processes)

  - For example, a **key-value record** managed by a storage server (or more complex data structure)

**Let's dissect the issue in detail**

# Shared data & synchronization



*What may happen if multiple threads concurrently access shared state (e.g., a shared variable)?*

# **Concurrently manipulating shared data**

- Two threads execute concurrently as part of the same process

- Shared variable (e.g., global variable)
  - `counter = 5`
- Thread 1 executes

  - `counter++`
- Thread 2 executes

  - `counter--`

- *What are **all the possible values** of* `counter` *after Thread 1 and Thread 2 executed?*

# Machine-level implementation

- Implementation of "`counter++`"

    **register$_1$ = counter**

    **register$_1$ = register$_1$ + 1**

    **counter = register$_1$**

- Implementation of "`counter--`"

    **register$_2$ = counter**

    **register$_2$ = register$_2$ – 1**

    **counter= register$_2$**

# Possible execution sequences

counter++

Context Switch

Context Switch

counter--

Context Switch

Context Switch

counter--

Context Switch

counter++

# Interleaved execution

- Assume **counter** is 5 and interleaved execution of counter++ (in $T_1$) and counter-- (in $T_2$)

  $T_1$: $r_1$ $\quad$ = **counter** $\quad$ (*register$_1$ = 5*)

  $T_1$: $r_1$ $\quad$ = $r_1$ + 1 $\quad$ (***register$_1$ = 6***)

  $T_2$ : $r_2$ $\quad$ = **counter** $\quad$ (*register$_2$ = 5*)

  $T_2$ : $r_2$ $\quad$ = $r_2$ – 1 $\quad$ (***register$_2$ = 4***)

  $T_1$ : **counter** $\quad$ = $r_1$ $\quad$ (***counter = 6***)

  $T_2$ : **counter** $\quad$ = $r_2$ $\quad$ (***counter = 4***)

  **Context switch**

- The value of **counter** may be 4 or 6, whereas **the correct result should be 5**!

# Race condition

- **Race condition**
  - **Several** threads **manipulate shared data concurrently**. The **final value** of the data **depends** upon **which thread finishes last**.

- In our example (interleaved execution) of `counter++` with `counter--`

- To prevent race conditions, concurrent processes must be **synchronized**!

# The moral of this story

- The statements

  **counter++;**

  **counter--;**

  must each be executed *atomically*.

- Atomic operation means an operation that **completes in its entirety without interruption**.

- This is achieved through **synchronization primitives**

- Shared variable accessed in a **critical section** must be protected by synchronization primitives

- Known as the **critical section problem** or as **mutual exclusion**

# Self-study questions

- Do we have a critical section problem in distributed systems? – There is no shared memory!

- Asked differently, do we need to worry about (distributed) mutual exclusion in a distributed system?

- Identify a few mutual exclusion scenarios in distributed systems.

"*Don't worry*", in distributed systems (the ones we look at), there is no shared memory – but,

...

# DISTRIBUTED MUTUAL EXCLUSION

## OVERVIEW

# Distributed mutual exclusion

- In distributed systems, mutual exclusion is at least equally complex due:
  - Lack of shared memory
  - Lack of a global clock
  - Event ordering

- Examples
  - Accessing a shared resource in distributed systems
  - Acquiring a lock
  - One active master to coordinate activities

# Critical section problem
## No shared memory

- System with *n* nodes

- Nodes access shared resources in CS

- **Coordinate access to CS via message passing**

- Application-level protocol for accessing CS

  – Enter_CS() – enter CS, block if necessary

  – ResourceAccess() – access shared resource in CS

  – Exit_CS() – leave CS

# Assumptions
## No practical rather theoretical considerations

- System is asynchronous
  - No bound on delays, no bound on clock drift, etc.
- Nodes do not fail
- Message delivery is reliable
  - Any message sent, is eventually delivered intact and exactly once – i.e., not lost, not duplicated
- Nodes are well-behaved and spent finite time accessing resources in CS

# Mutual exclusion requirements

- **Safety** – correctness
  - At most one node in the critical section at a time
- **Liveness –** progress (something good happens)
  - Requests to enter/exit CS eventually succeed
  - No deadlock
- **Fairness** (order & starvation)
  - If one request to enter CS **happened-before** another one, then entry to CS is granted in that **order**
  - Requests are ordered such that no node enters the critical section twice while another waits to enter (i.e., **no starvation**)

# Deadlock & starvation

- **Deadlock**: Two or more nodes become **stuck indefinitely** while attempting to enter and exit CS – e.g., by virtue of their mutual dependency

- **Starvation**: The **indefinite postponement** of entry to CS for a node that has requested it.

$R_1$

$R_2$

Holds resource

Holds resource

$N_1$

$N_2$

Mutually wait for resources held by another node

# Possible performance metrics

- **Bandwidth:** Number of messages sent, received or both

- **Synchronization delay**: Time between one process **exiting** critical section and next one **entering**

- **Client delay**: Delay **at entry** and **exit** (response time)

- We do not measure client access to resources protected by the critical section (**assume finite**)

# Solution strategy overview

- **Centralized strategy**
  - Divide nodes into **leader** and **follower**, leader dictates actions of followers

- **Distributed strategy**: Each node independently decides actions, based on local knowledge of others' state
  - **Token-based**: A node is allowed in the critical section (CS) if it has a token. Tokens are passed from node to node, in some (priority) order.
  - **Non-token-based**: A node enters CS when an **assertion becomes true**. A node communicates with other nodes to obtain their states and decides whether the assertion is true or false.

# Self-study questions

- What are some examples for mutual exclusion scenarios in distributed systems?

- What are some scenarios where we need locks in distributed systems, after all, there is no shared memory (in our notions of DS)

Pixabay license

# DISTRIBUTED MUTUAL EXCLUSION

## CENTRALIZED STRATEGY

# Centralized strategy: Elect leader

Elect a leader (not in scope of this strategy)

# Centralized strategy: **Elect leader**

Elect a leader (not in scope of this strategy)

# Centralized strategy: Empty CS

A.k.a. server / coordinator / **leader** / master

$P_5$

**Critical section**

$P_1$    $P_2$  ...  $P_i$  ...  $P_n$

# Centralized strategy: Empty CS

A.k.a. server / coordinator / **leader** / master



P₅

**Critical section**

Request entry to CS
**Enter_CS()**

P₁   P₂   ...   Pᵢ   ...   Pₙ

# Centralized strategy: Empty CS

A.k.a. server / coordinator / **leader** / master

$P_5$

**Critical section**

$P_1$  $P_2$  ...  $P_i$  ...  $P_n$

# Centralized strategy: **Empty CS**

A.k.a. server / coordinator / **leader** / master



$P_5$

**Critical section**

Grant access to CS
(pass the **token**)

$P_1$  $P_2$  ...  $P_i$  ...  $P_n$

# Centralized strategy: Empty CS

A.k.a. server / coordinator / **leader** / master

$P_5$

**Critical section**

$P_1$   $P_2$   ...   $P_i$   ...   $P_n$

# Centralized strategy: Empty CS

A.k.a. server / coordinator / **leader** / master

**$P_2$**

**Critical section**

$P_5$

$P_1$    $P_2$  ...  $P_i$  ...  $P_n$

**ResourceAccess()**

# Centralized strategy: Empty CS

A.k.a. server / coordinator / **leader** / master

# Centralized strategy: Empty CS

A.k.a. server / coordinator / **leader** / master

$P_5$

**Critical section**

$P_1$   $P_2$   ...   $P_i$   ...   $P_n$

# Centralized strategy: Non-empty CS

# Centralized strategy: **Non-empty CS**

# Centralized strategy: Non-empty CS

# Centralized strategy: **Non-empty CS**



$P_2$

**Critical section**

Queue $P_1$ $P_n$

$P_5$

Request entry to CS
**Enter_CS()**

$P_1$ $P_2$ ... $P_i$ ... $P_n$

**ResourceAccess()**

# Centralized strategy: **Non-empty CS**

# Centralized strategy: Exit CS



Queue    P₁  Pₙ    P₅

P₂

**Critical section**

P₁    P₂  ...  Pᵢ  ...  Pₙ

**ResourceAccess()**

# Centralized strategy: **Exit CS**



Critical section

Queue  P$_1$  P$_n$

P$_5$

Exit_CS()

P$_1$  P$_2$  ...  P$_i$  ...  P$_n$

# Centralized strategy: **Exit CS**



Critical section

Queue $P_1$ $P_n$

$P_5$

$P_1$    $P_2$    ...    $P_i$    ...    $P_n$

# Centralized strategy: **Exit CS**



Critical section

P₅

Queue    P₁    Pₙ

Grant access to CS
(pass the **token**)

P₁    P₂    ...    Pᵢ    ...    Pₙ

# Centralized strategy: Exit CS



Queue $P_1$   $P_5$

$P_n$
**Critical section**

Grant access to CS
(pass the **token**)

$P_1$   $P_2$   ...   $P_i$   ...   $P_n$

**ResourceAccess**

# Summarizing observations I

- Meets requirements: Safety, liveness, no starvation
- ***Does solution meet the ordering requirement?***
- Advantages
  - **Simple to implement**
- Disadvantages
  - **Single point of failure**
  - Bottleneck, network congestion, timeout
- Deadlock potential for multiple resources with separate servers

# Summarizing observations II

- Enter_CS()
  - **Two messages**: Request & Grant
  - One round of communication (RTT delay)
- Exit _CS()
  - **One message**: Release message
  - No delay for the node in CS

# Self-study questions

- Does solution meet the ordering requirement? Why or why not?

- How could a single leader manage multiple resources? Analyse pros and cons of alternative designs?

- Provide a differentiated discussion regarding failure of leader, follower, follower in CS, follower requested access, etc.)

Token-ring

# DISTRIBUTED MUTUAL EXCLUSION
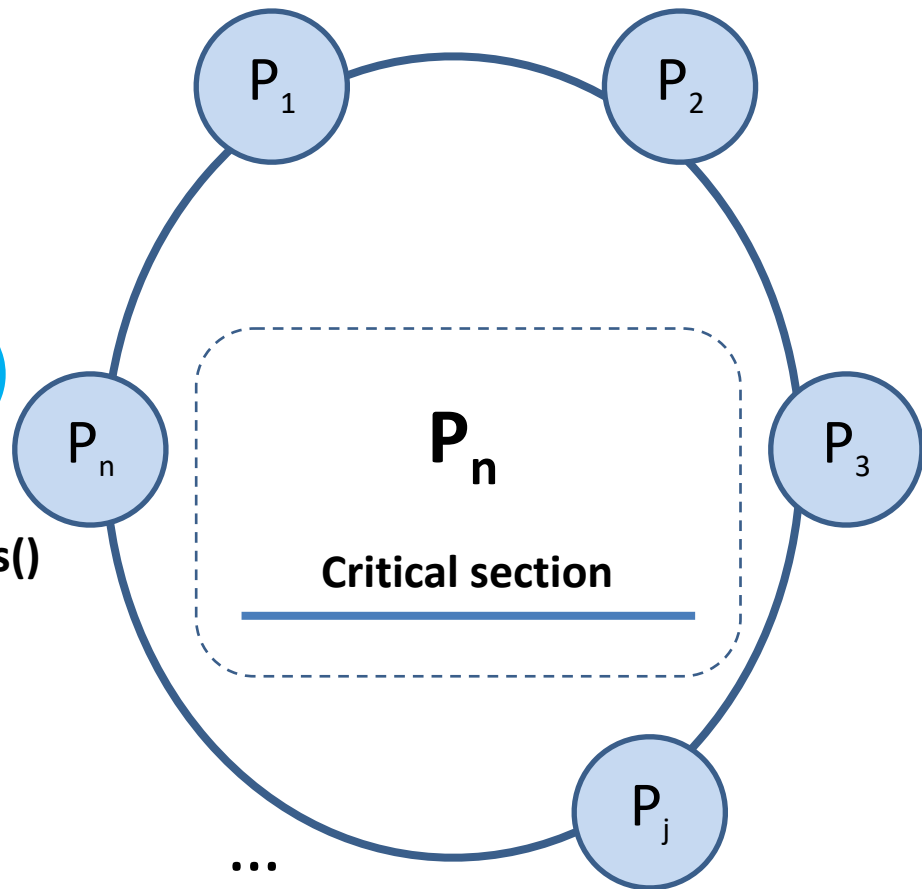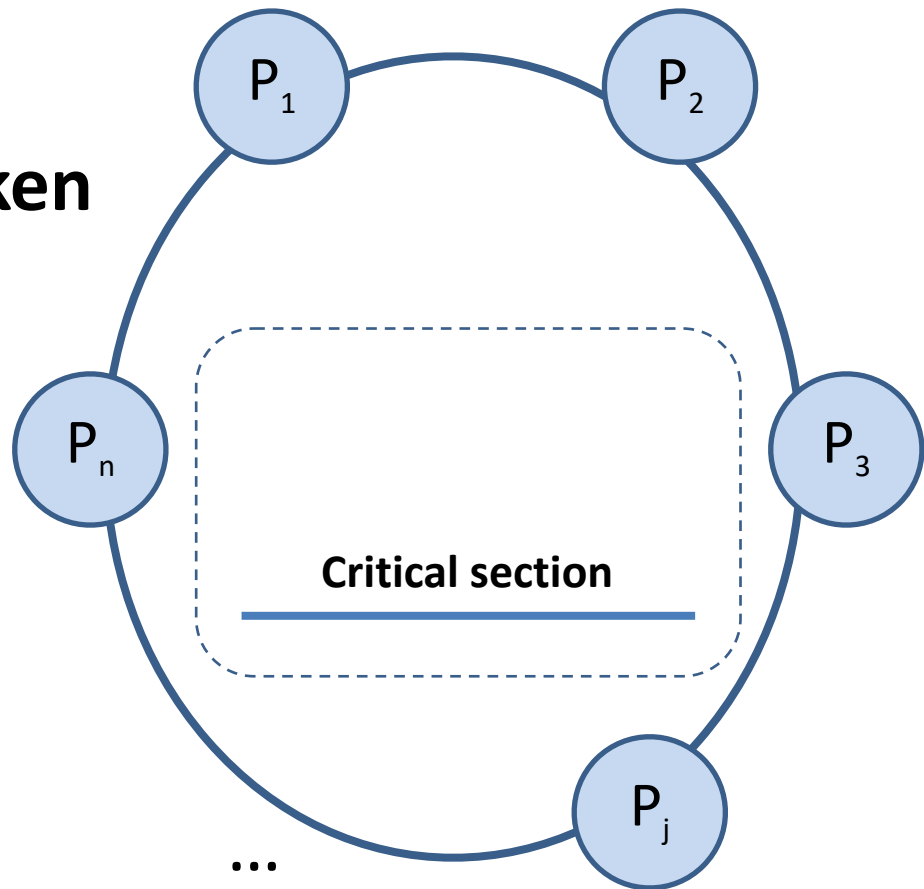
## RING-BASED ALGORITHM

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \mod N}$
- Logical topology a priori unrelated to physical topology



$P_1$  $P_2$  $P_3$  $P_j$  $P_n$

**Critical section**

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology



**Token**

$P_1$   $P_2$   $P_3$   $P_j$   $P_n$

**Critical section**

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Enter_CS()**

**Token**

$P_1$

$P_2$

$P_n$

$P_3$

$P_j$

**Critical section**

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Enter_CS()**

P$_1$  P$_2$  P$_n$  P$_3$  P$_j$

**Critical section**

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its
  successor, $P_{(i+1) \bmod N}$
- Logical topology
  a priori unrelated
  to physical **Enter_CS()**
  topology



$P_1$

$P_2$

$P_n$

$P_3$

**Critical section**

$P_j$

...

**Token**

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Enter_CS()**



$P_1$  $P_2$  $P_n$  $P_3$  $P_j$

**Critical section**

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Enter_CS()**



P₁  P₂  Pₙ  P₃  Pⱼ

**Critical section**

**Token**

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Enter_CS()**

$P_1$   $P_2$

$P_n$   $P_3$

**Critical section**

$P_j$

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Enter_CS()**

**Token**



P_1, P_2, P_3, P_j, P_n

Critical section

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Enter_CS()**

**Token**

ResourceAccess()



$P_1$  $P_2$

$P_n$  $P_3$

$P_j$

$P_n$

**Critical section**

...

# Ring-based algorithm

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

**Token**

**Exit_CS()**

$P_1$ $P_2$ $P_n$ $P_3$

**Critical section**

$P_j$

...

# Ring-based algorithm analysis

- **Safe**: Node enters CS only if it holds token

- **Live**: Since finite work is done by each node (can't re-enter), token eventually gets to each node

- **Fair**: Ordering is based on ring topology, no starvation (pass token between accesses)

- **Performance**

  — **Constantly consumes network bandwidth**, even when no node seeks entry, except when inside CS

  — **Synchronization delay**: Between 1 and $N$ messages

  — **Client delay**: 0 to $N$ messages for entry; 0 for exit

# Potential problems with ring-based algorithm
## (Due to our assumption, not all apply here.)

- Node crash

- Lost token

- Duplicate token

- Timeouts on token passing

# Self-study questions

- Does solution meet the ordering requirement? Why or why not?

- How could access to multiple resources be realized with the ring-based ME algorithm?

- How can node failures be mitigated? Develop strategies that could tolerate 1, 2, … node failures?

- How could lost or duplicated tokens be mitigated?

# DISTRIBUTED MUTUAL EXCLUSION

## LAMPORT'S ALGORITHM FOR MUTUAL EXCLUSION, 1978

# Lamport's algorithm
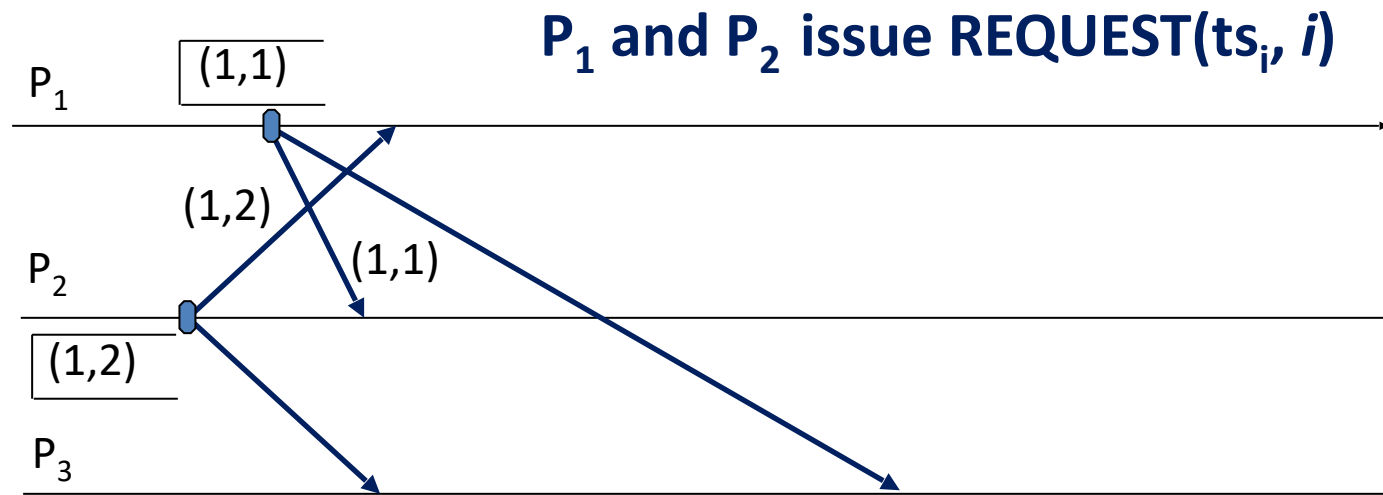


Queues likely
Pixabay license

- System of $n$ nodes

- $(ts_i, i)$ – logical clock timestamp of node $P_i$ (**total order!**)

- Logical timestamps serve to order requests for CS

- Smaller timestamps have priority over larger ones

- Request queue maintained at each node ordered by timestamp (**a priority queue!**)

- Assume message delivered in FIFO order, i.e., messages do not race over communication channel
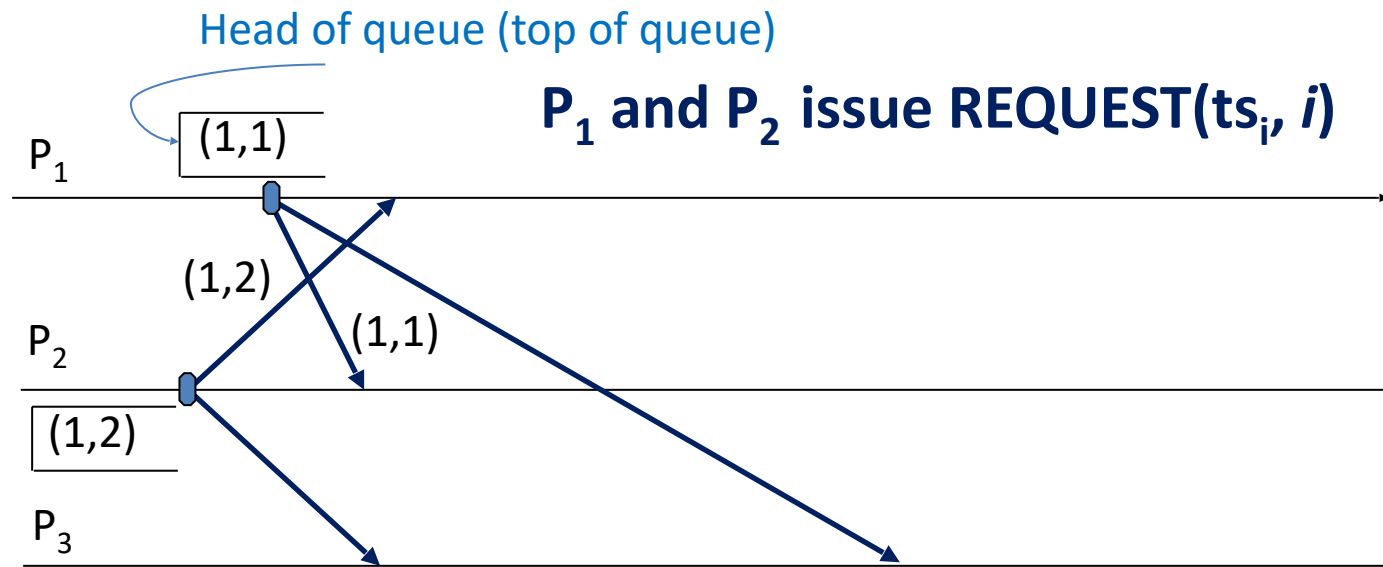
# $P_i$  Lamport's algorithm: Request

N.B.: $(ts_i, i)$ denotes timestamp of request

**$P_i$** requesting CS

— Broadcast **REQUEST($ts_i$, *i*) message** to *all* nodes

  • Place request in its ***request_queue$_i$***

**$P_1$ and $P_2$ issue REQUEST($ts_i$, *i*)**



$P_1$  (1,1)

(1,2)

(1,1)

$P_2$

(1,2)

$P_3$

# $\mathbf{P_i}$   **Lamport's algorithm: Request**

N.B.: ($ts_i$, $i$) denotes timestamp of request

**$P_i$ requesting CS**

— Broadcast **REQUEST($ts_i$, $i$) message** to **all** nodes
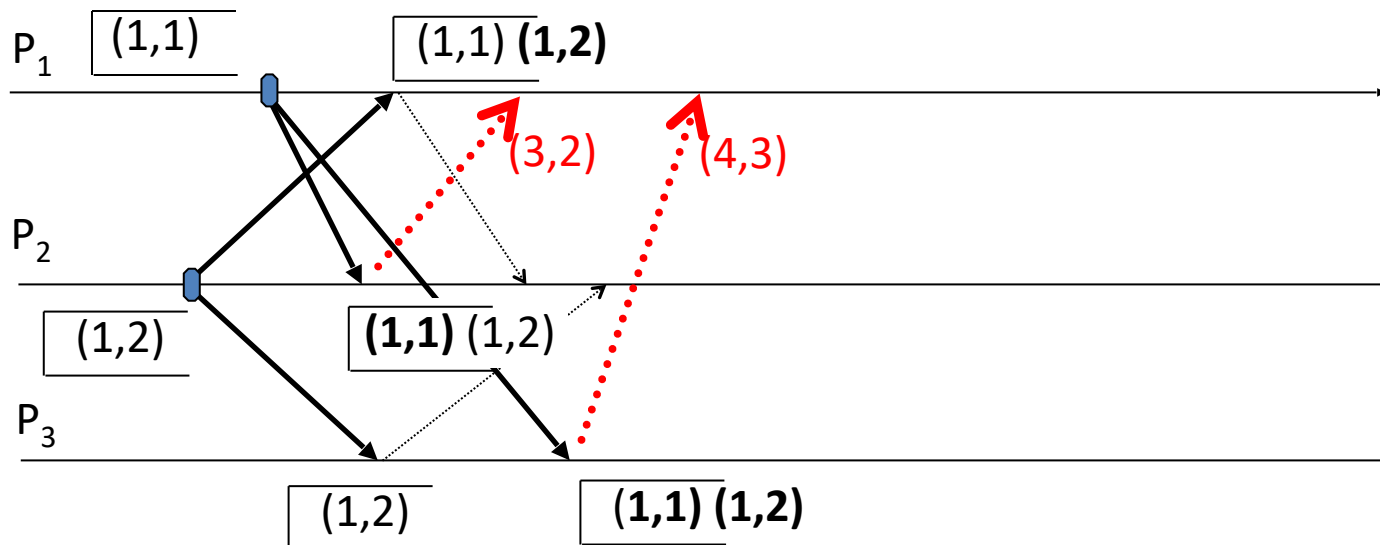
  • Place request in its ***request_queue$_i$***

Head of queue (top of queue)

**$P_1$ and $P_2$ issue REQUEST($ts_i$, $i$)**

(1,1)

$P_1$

(1,2)

(1,1)

$P_2$

(1,2)

$P_3$

# Lamport's algorithm: Receive request

$P_k$

**$P_k$** receiving a request to enter CS

- When **$P_k$** receives a **REQUEST($ts_i$, *i*) message** from **$P_i$**
  - It places **$P_i$**'s request into its ***request_queue$_k$***
  - Sends a timestamped **REPLY message** to **$P_i$**

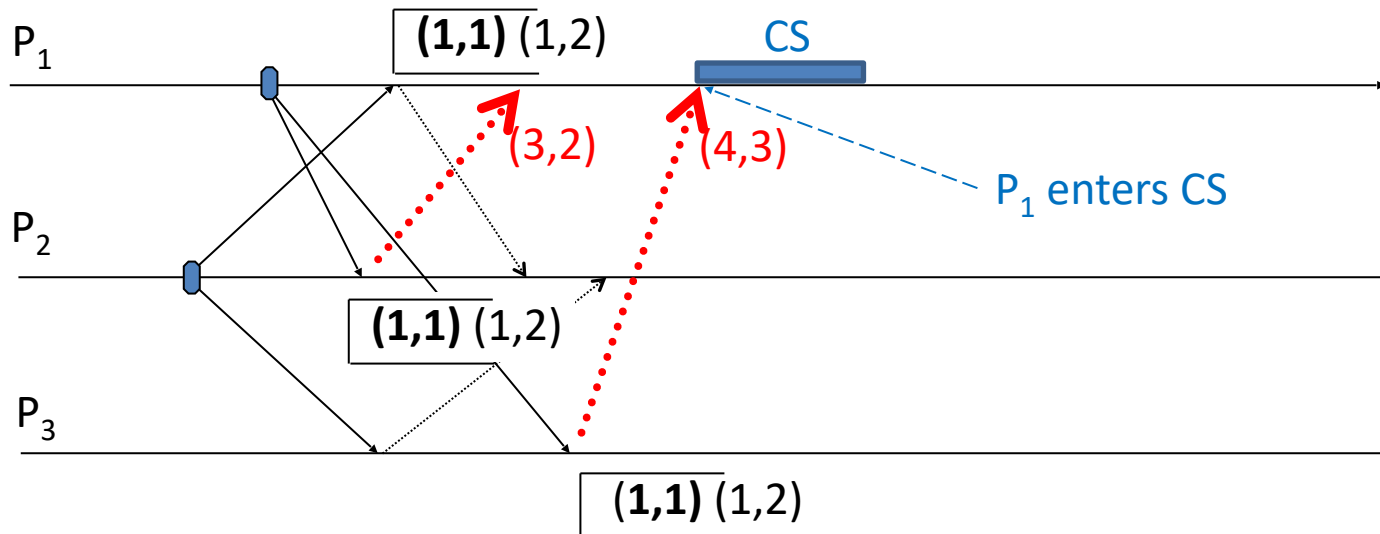**REPLY($ts_k$, *k*) – only timestamped replies to $P_i$'s request are shown below**

# $P_i$ Lamport's algorithm

$P_i$ requesting CS – **$P_i$ enters CS when following conditions hold:**

1. **$P_i$ received a message with timestamp larger than $(ts_i, i)$ from every other node**

2. **$P_i$'s request is at top of *request_queue_i***

**$P_1$ sent REQUEST(1, 1);** **now, it received REPLY(3,2) and REPLY(4,3), and its request (1, 1) is at the top of queue**

Its original request timestamp



CS

$P_1$ enters CS

# $P_i$ Lamport's algorithm

**$P_i$ releasing CS**

— Removes its request from top of the request queue

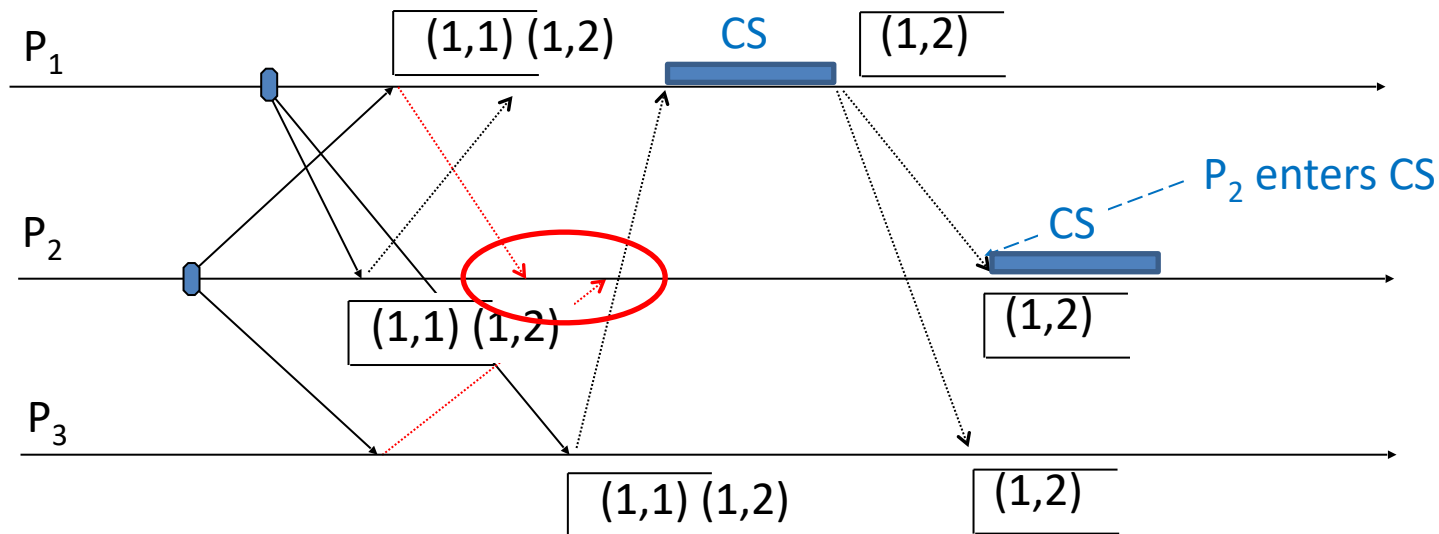— Broadcasts a timestamped **RELEASE** message to all nodes

**$P_1$ - RELEASE(ts$_i$, *i*)**

# $P_k$ Lamport's algorithm

**$P_k$** receiving a release message

- When $P_k$ receives a **RELEASE** message from $P_i$, $P_k$ removes $P_i$'s request from its request queue

# **Summarizing observations**

- 3(N-1) messages per CS request invocation
    - (N - 1) REQUEST
    - (N - 1) REPLY
    - (N - 1) RELEASE messages

- Not fault tolerant

# Self-study questions

- Would node IDs as timestamps suffice, why or why not?

- Would Lamport clock values as timestamps suffice, why or why not?

- Is the algorithm order-preserving, why or why not?

# DISTRIBUTED MUTUAL EXCLUSION

## RICART & AGRAWALA, 1981

# Ricart & Agrawala, 1981, algorithm
## (Guarantees mutual exclusion among *n* node)

- **Basic idea**
  - Nodes wanting to enter CS, **broadcast a request to all other nodes**
  - Enter CS, once **all** nodes have **granted request**
- Use **Lamport timestamps** to order requests: $<ts_i, i>$, $ts_i$ the timestamp, $i$ the node identifier

# Ricart & Agrawala, 1981, algorithm
## (Guarantees mutual exclusion among *n* node)

- **Basic idea**
  - Nodes wanting to enter CS, **broadcast a request to all other nodes**
  - Enter CS, once **all** nodes have **granted request**
- Use **Lamport timestamps** to order requests: $<ts_i, i>$, $ts_i$ the timestamp, $i$ the node identifier

# Ricart & Agrawala, 1981, algorithm
## (Guarantees mutual exclusion among *n* node)

- **Basic idea**
  - Nodes wanting to enter CS, **broadcast a request to all other nodes**
  - Enter CS, once **all** nodes have **granted request**
- Use **Lamport timestamps** to order requests: $<ts_i, i>$, $ts_i$ the timestamp, $i$ the node identifier

# Ricart & Agrawala, 1981, algorithm
## (Guarantees mutual exclusion among *n* node)

- **Basic idea**
  - Nodes wanting to enter CS, **broadcast a request to all other nodes**
  - Enter CS, once **all** nodes have **granted request**
- Use **Lamport timestamps** to order requests: $<ts_i, i>$, $ts_i$ the timestamp, *i* the node identifier

# Ricart & Agrawala: Distributed strategy

- Each node is in one of three states
  - **Released**      - Outside the CS, e.g., after Exit_CS()
  - **Wanted**       - Wanting to enter CS, in call Enter_CS()
  - **Held**          - Inside CS, during RessourceAccess()

- If a node requests to enter CS and **all** other **nodes** are **in the *Released state*, entry** is **granted** by each node

- If a node, $P_i$, requests to enter CS and another node, $P_k$, is inside the CS (***Held state***), then $P_k$ will not reply, until it is finished with the CS

# Initialization

**RELEASED**

**Critical section**

$P_1$

...

$P_2$

**RELEASED**

$P_3$

**RELEASED**

# Requesting entry to CS

## Request while *all Released*

**Critical section**

**RELEASED**

P₁

...

P₂

**WANTED**

P₃

**RELEASED**

# Requesting entry to CS

## Request while *all Released*



RELEASED

Critical section

P₁

...

Request

Request

P₂

P₃

WANTED

RELEASED

# Requesting entry to CS

## Request while *all Released*

Critical section

RELEASED

P₁

...

WANTED

RELEASED

# Requesting entry to CS

## Request while *all Released*

**RELEASED**

**Critical section**

$P_1$

...

Reply

$P_2$

$P_3$

**WANTED**

**RELEASED**

# Requesting entry to CS

## Request while *all Released*

**Critical section**

**RELEASED**

$P_1$

...

$P_2$

**WANTED**

$P_3$

**RELEASED**

# Requesting entry to CS

## Request while *all Released*

**Critical section**

**RELEASED**

$P_1$

...

$P_2$ ← Reply — $P_3$

**WANTED**                    **RELEASED**

# Requesting entry to CS

## Request while *all Released*

**Critical section**

**RELEASED**

P<sub>1</sub>

...

P<sub>2</sub>

**WANTED**

P<sub>3</sub>

**RELEASED**

# Requesting entry to CS

## Request while *all Released*

$P_2$

**Critical section**

**RELEASED**

$P_1$

...

$P_2$

**HELD**

$P_3$

**RELEASED**

# Requesting entry to CS

## *Request while Held*



P₃

Critical section

RELEASED

P₁

...

P₂

**WANTED**

P₃

**HELD**

# Requesting entry to CS

## *Request while Held*

# Requesting entry to CS
## *Request while Held*



P₃

**Critical section**

**RELEASED**

P₁

...

P₂

**WANTED**

P₃

**HELD**

# Requesting entry to CS

## *Request while Held*



RELEASED

P₃

**Critical section**

P₁

...

Reply

P₂

WANTED

P₃

HELD

# Requesting entry to CS

## Request while Held

# Requesting entry to CS

## Request while Held



Critical section

RELEASED

$P_1$

...

Reply

Exit_CS()

$P_2$

$P_3$

Reply

WANTED

RELEASED

# Requesting entry to CS

## *Request while Held*

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry**
**to CS concurrently**

**RELEASED**

...

P$_1$

P$_2$

**WANTED**

P$_3$

**RELEASED**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry**
**to CS concurrently**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry to CS concurrently**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry
to CS concurrently**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry
to CS concurrently**

RELEASED

...  P$_1$

P$_2$

**WANTED**

P$_3$

**WANTED**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry
to CS concurrently**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry to CS concurrently**

RELEASED

...

P$_1$

P$_2$

**WANTED**

P$_3$

**WANTED**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry**
**to CS concurrently**

RELEASED

P$_1$

**@P$_3$**: 7 from **P$_3$** < 15 from P$_2$,

Own timestamp lower, therefore,

hold on to reply to P$_2$

<7, 3> < <15, 2>

P$_2$

P$_3$

Queue

Req.: <15, 2>

WANTED

WANTED

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry
to CS concurrently**

RELEASED

P$_1$

**@P$_2$**: 7 from P$_3$ < 15 from **P$_2$**,
therefore grant P$_3$
access first

<7, 3> < <15, 2>

**@P$_3$**: 7 from **P$_3$** < 15 from P$_2$,
Own timestamp lower, therefore,
hold on to reply to P$_2$

P$_2$

WANTED

P$_3$

Queue

Req.:  <15, 2>

WANTED

# *Concurrent* entry requests

**P₂ and P₃ request entry**
**to CS concurrently**

**RELEASED**

**@P₂**: 7 from P₃ < 15 from **P₂**, therefore grant P₃ access first

**@P₃**: 7 from **P₃** < 15 from P₂, Own timestamp lower, therefore, hold on to reply to P₂



<7, 3> < <15, 2>

P₁

P₂

P₃

Reply

**WANTED**

**WANTED**

Queue

Req.: <15, 2>

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry**

**to CS concurrently**

RELEASED

P$_1$

@**P$_3$**: 7 from **P$_3$** < 15 from P$_{2,}$

Own timestamp lower, therefore,

hold on to reply to P$_2$

<7, 3> < <15, 2>

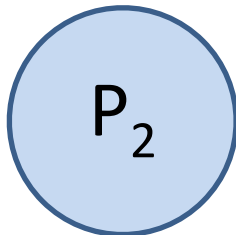P$_2$

P$_3$

Queue

Req.:  <15, 2>

**WANTED**

**WANTED**

# *Concurrent* entry requests

**P₂ and P₃ request entry
to CS concurrently**

RELEASED

P₃

Critical section

$P_1$

**@P₃**: 7 from **P₃** < 15 from P₂,
Own timestamp lower, therefore,
hold on to reply to P₂

<7, 3> < <15, 2>

$P_2$

$P_3$

Queue

Req.: <15, 2>

**WANTED**

**HELD**

# *Concurrent* entry requests

**P<sub>2</sub> and P<sub>3</sub> request entry to CS concurrently**

RELEASED

$P_1$

$P_3$

**Critical section**

$P_2$

**WANTED**

$P_3$

**HELD**

Queue

Req.:  <15, 2>

# *Concurrent* entry requests

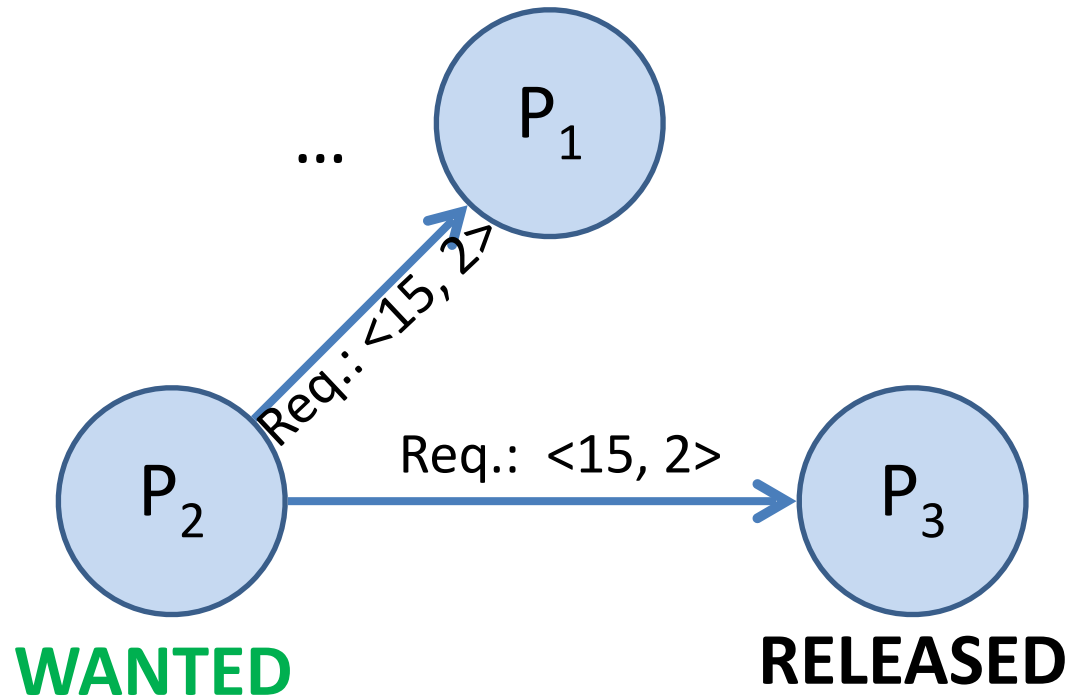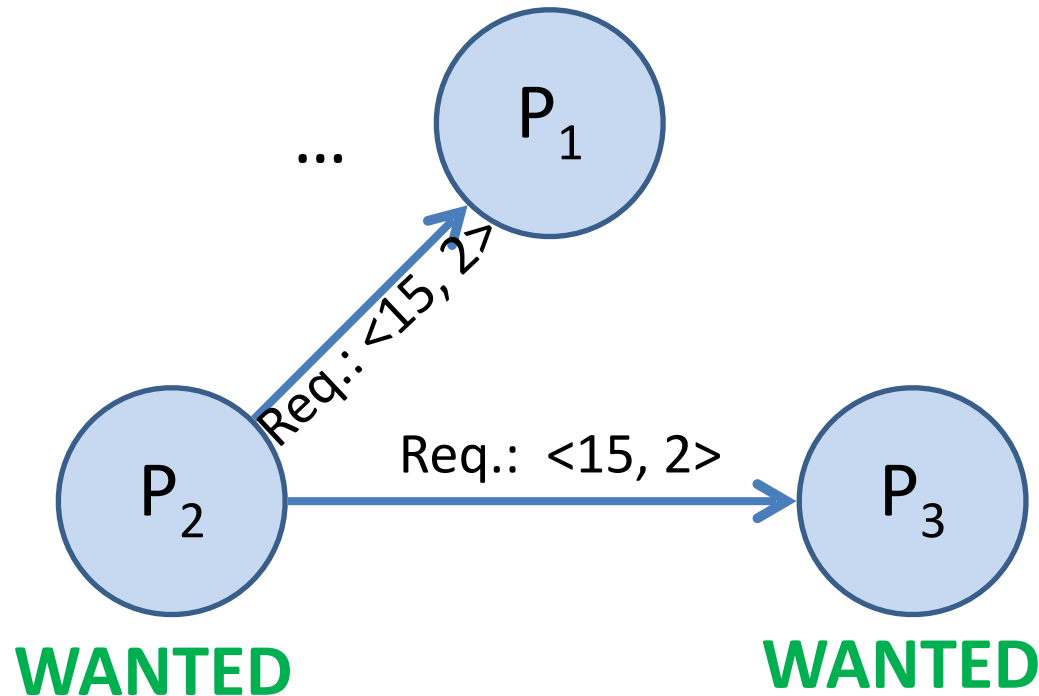**P$_2$ and P$_3$ request entry to CS concurrently**

# *Concurrent* entry requests
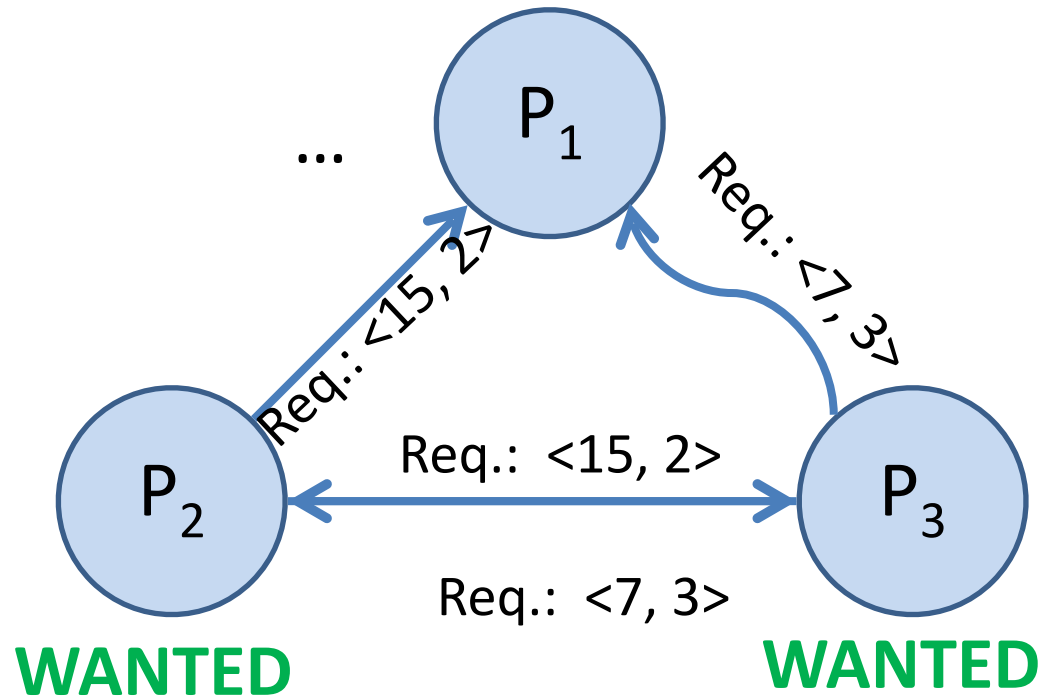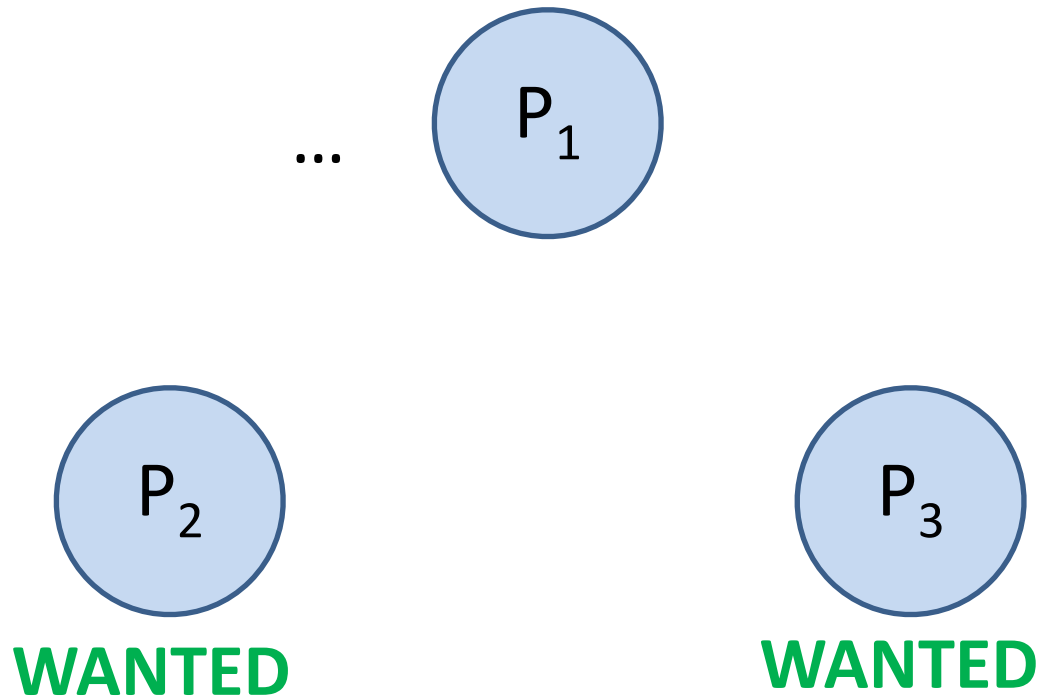
**P$_2$ and P$_3$ request entry to CS concurrently**
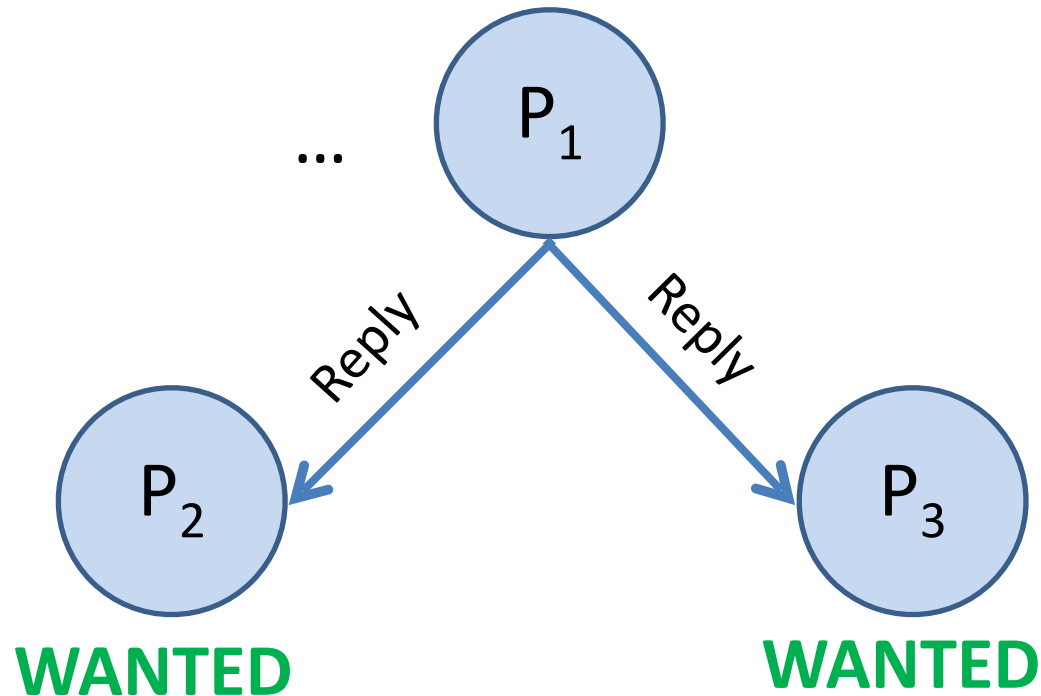
RELEASED

P$_1$

Critical section

P$_2$

**WANTED**

P$_3$

**RELEASED**

# *Concurrent* entry requests

**P$_2$ and P$_3$ request entry to CS concurrently**

**RELEASED**

P$_1$

P$_2$

**Critical section**

P$_2$

**HELD**

P$_3$

**RELEASED**

# Pseudo code

# Pseudo code

**On initialization**
   state = RELEASED

# Pseudo code

**Enter_CS()**

    state = WANTED

    Broadcast timestamped
        **request** to all nodes

    wait until ( (n-1) acks received )

    state = HELD

**On initialization**

    state = RELEASED

# Pseudo code

**Enter_CS()**

    state = WANTED

    Broadcast timestamped
        **request** to all nodes

    wait until ( (n-1) acks received )

    state = HELD

**On initialization**

    state = RELEASED

**On receiving a request with $<ts_i, i>$, at $P_k$ (i ≠ k)**

    if ( state==HELD or ( state==WANTED and $<ts_k, k>$ < $<ts_i, i>$ ) )

        queue request from $P_i$ without replying

    else

        send a reply to $P_i$

**k**

# Pseudo code

**Enter_CS()**

    state = WANTED

    Broadcast timestamped

        **request** to all nodes

    wait until ( (n-1) acks received )

    state = HELD

**On initialization**

    state = RELEASED

**Exit_CS()**

    state = RELEASED

    Reply to all queued

    requests

**On receiving a request with $<ts_i, i>$, at $P_k$ (i ≠ k)**

    if ( state==HELD or ( state==WANTED and $<ts_k, k> < <ts_i, i>$ ) )

        queue request from $P_i$ without replying

    else

        send a reply to $P_i$

**k**

# Pseudo code

**Enter_CS()**

    state = WANTED

    Broadcast timestamped

        **request** to all nodes

    wait until ( (n-1) acks received )

    state = HELD

**On initialization**

    state = RELEASED

**Exit_CS()**

    state = RELEASED

    Reply to all queued

        requests

*k*'s timestamp
(its own)

**On receiving a request with $<ts_i, i>$, at $P_k$ (i ≠ k)**

    if ( state==HELD or ( state==WANTED and $<ts_k, k> < <ts_i, i>$ ) )

        queue request from $P_i$ without replying

    else

        send a reply to $P_i$

**k**

# Pseudo code

**Enter_CS()**

    state = WANTED

    Broadcast timestamped
        **request** to all nodes

    wait until ( (n-1) acks received )

    state = HELD

**On initialization**

    state = RELEASED

**Exit_CS()**

    state = RELEASED

    Reply to all queued
    requests

**On receiving a request with $<ts_i,\ i>$, at $P_k$ (i ≠ k)**

    if ( state==HELD or ( state==WANTED and $<ts_k,\ k> < <ts_i,\ i>$ ) )

        queue request from $P_i$ without replying

    else

        send a reply to $P_i$

**k**

# Reminder: Subtlety about timestamps

- Use **Lamport timestamps** to order requests: $<ts_i, i>$, $ts_i$ the timestamp, $i$ the node identifier

- If for two timestamps
  - $<ts_i, i>$ and $<ts_k, k>$, if $ts_i = ts_i$, then break ties by looking at node identifiers $i, k$
  - Gives rise to an **arbitrary total order** over timestamps (i.e., requests)

# Summarizing observations I

- **Safe**
  - If two nodes were in the CS at the same time, they would have had to reply to each other, which can't happen since their requests are totally ordered

- **Live**
  - Each request is enqueued, so eventually it will be satisfied

- **Ordered**
  - Based on timestamps

# Summarizing observations II

- Each entry request requires 2(N-1) messages
  - N-1 requests, N-1 replies
- Synchronization delay is one message
- Not fault tolerant

# Self-study questions

- Would node IDs as timestamps suffice, why or why not?

- Would Lamport clock values as timestamps suffice, why or why not?

- Is the algorithm order-preserving, why or why not?

- Compare all our distributed mutual exclusion algorithms according to bandwidth use, synchronization delay, and delay at entry and exit.

- Compare all our distributed mutual exclusion algorithms according to their failure resilience.

Pixabay license

# LEADER ELECTION
## OVERVIEW

# Leader election
## a.k.a., coordinator, master, etc.

- **Problem**: A group of nodes, $P_1$, …, $P_n$, **must agree** on some **unique $P_k$** to be the "**leader**"

- Often, leader then **coordinates another activity**

- Election runs when leader failure is detected (suspected)

- Any node who hasn't heard from leader in some predefined time interval **may call for an election**

- **False alarm** is a possibility (new elections initiated, while current leader still alive)

- **Several nodes** may initiate **elections concurrently**

- Algorithm should allow for node crash during election

# **Leader election use cases**

- Berkeley clock synchronization algorithm

- Centralized mutual exclusion algorithm

- Leader election for choosing *master* in Hbase, Bigtable

- Choosing *master* among *n* nodes in Chubby or ZooKeeper coordination service


- *Primary-backup replication algorithms*

- *Two-phase commit protocol*

# Leader election vs. mutual exclusion

- Election losers return to what they were doing …

  … **instead of waiting**

- **Fast election** is important …

  … not starvation avoidance

- **All nodes** must know result …

  … not just the winner

# Leader election vs. mutual exclusion

- Election losers return to what they were doing …

    … **instead of waiting**

- **Fast election** is important …

    … not starvation avoidance

- **All nodes** must know result …

    … not just the winner

> ME can be reduced to LE!
> (e.g., HBase wants LE,
> ZooKeeper provides ME)

# Uniqueness requirement
## Unique identifier

- Elected **leader** must be **unique**
- Active nodes with **largest identifier** wins
- Unique identifier (UID) could be any "useful value"
  - I.e., **unique** and **totally orderable value**
- E.g., based on process identifiers, IP-port-number
- E.g., based on least computational load
  - $<1/load, i>$, $load > 0$, $i$ is UID to break ties
- Each node, $P_i$, has a variable $elected_i$ that holds the value of the leader or is "$\perp$" (undefined)

# Election algorithm requirement

- **Safety** - correctness

  - A participating node $P_i$ has variable $elected_i$ = " $\perp$ "  or $elected_i$ = $P$, where $P$ is chosen as the non-crashed node at the end of the election run with the **largest identifier.**

  - **Only one leader at a time!**

- **Liveness** – progress is made

  - **All nodes participate** in the election and **eventually** either set $elected_i$ ≠ " $\perp$ " or **crash**.

# Summary

- Leader election is a fundamental building block in designing distributed systems

- Leader election algorithms overview (a popular shortlist)
  - Chang & Roberts, 1979
  - HS algorithm, 1980
  - Bully, 1982
  - Leader Election in Raft, 2014

# Self-study questions

- Why settle on the largest unique identifier as determining characteristic for leader?

- Develop your own leader election algorithm by reducing the problem to mutual exclusion.

- Find a few more leader election use case scenarios.

Pixabay license

# LEADER ELECTION
## CHANG & ROBERTS – RING-BASED ALGORITHM, 1979

# Assumptions and setup

- Construct a ring (cf. ring-based mutual exclusion)
- Assume, each *P* has **unique identifier (UID)**
- Assume, no failures and asynchronous system, **but failures do happen before the election!**
- Any $P_i$ may begin an election by sending an **election message** to its successor (i.e., after suspecting a leader failure)
- Election message holds $P_i$'s UID

**Ring-based algorithm**

- Logical ring of nodes
- Each $P_i$ knows its successor, $P_{(i+1) \bmod N}$
- Logical topology a priori unrelated to physical topology

Distributed Systems (Hans-Arno Jacobsen)

# Ring-based election algorithm

- Election message holds $P_i$'s $UID_i$

- Nodes manage a flag **PARTICIPANT** to remember their election participation status (initially false)

- Upon receiving an election message, $P_k$ **compares its own $UID_k$ to $UID_i$**

  - If $UID_i$ > $UID_k$: **Forwards** election message

  - If $UID_i$ < $UID_k$: **Forwards** election message **with $P_k$'s $UID_k$** unless $P_k$ has already sent a message, i.e., has participated in current election run (in that case, it does nothing)

  - If $UID_i$ = $UID_k$: $P_k$ **is now leader**, forwards *victory message* to notify all other nodes (resets PARTICIPANT flag)

# Ring-based algorithm
## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (determine winner)

Non-participant

**Participant**

Non-participant

Non-participant

Non-participant

$P_1$

$P_2$

$P_n$

$P_3$

...

$P_j$

# Ring-based algorithm
## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (determine winner)

# Ring-based algorithm

## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (determine winner)



Non-participant

**Participant**

$P_1$

$P_2$

**2**

Non-participant

**Participant**

$P_n$

$P_3$

**3**

...

**Participant**

$P_j$

# Ring-based algorithm

## Calling an election (determine winner)



Non-participant

Non-participant

**Participant**

**Participant**

**Participant**

$P_1$

$P_2$

$P_3$

$P_n$

$P_j$

**2**

**3**

...

i

# Ring-based algorithm

## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (determine winner)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

# Ring-based algorithm
## Calling an election (origin & victory)

Elected n

Non-participant

Elected n

Non-participant

Leader

Participant

n

Participant

n

Participant

n

$P_1$ $P_2$ $P_3$ $P_j$ $P_n$

...

# Ring-based algorithm
## Calling an election (<span style="color:red">origin & victory</span>)

# Ring-based algorithm
## Calling an election (origin & victory)



Elected n

Non-participant

Non-participant

Elected n

Elected n

P₁

P₂

Non-participant

P₃

Non-participant

Leader

n

Elected n

Pₙ

=

Pⱼ

Participant

...

n

# Ring-based algorithm

## Calling an election (origin & victory)

S$_{ID}$  ID in message from sender

R$_{ID}$  ID at receiver

# Different cases

**S (sender)**　　　　　　**R (receiver)**

**S$_{ID}$ > R$_{ID}$**

**Participant**

S$_{ID}$ ID in message from sender

R$_{ID}$ ID at receiver

# Different cases

**S (sender)**         **R (receiver)**

$$\mathbf{S_{ID} > R_{ID}}$$

**Participant**
Forward S$_{ID}$

$S_{ID}$  ID in message from sender

$R_{ID}$  ID at receiver

# Different cases

**S (sender)**                    **R (receiver)**

$S_{ID} > R_{ID}$

**Participant**
 Forward $S_{ID}$

$S_{ID} > R_{ID}$

**Non-participant**

$S_{ID}$ ID in message from sender
$R_{ID}$ ID at receiver

# Different cases

**S (sender)**                    **R (receiver)**

$$S_{ID} > R_{ID}$$

**Participant**
  Forward $S_{ID}$

$$S_{ID} > R_{ID}$$

**Non-participant**
  Forward $S_{ID}$    $\rightarrow$ Participant

$S_{ID}$ ID in message from sender

$R_{ID}$ ID at receiver

# Different cases

**S (sender)**                    **R (receiver)**

$S_{ID} > R_{ID}$

**Participant**
Forward $S_{ID}$

$S_{ID} > R_{ID}$

**Non-participant**
Forward $S_{ID}$      $\rightarrow$ Participant

$S_{ID} < R_{ID}$

**Non-participant**

$S_{ID}$  ID in message from sender

$R_{ID}$  ID at receiver

# Different cases

**S (sender)**            **R (receiver)**

$\mathbf{S_{ID} > R_{ID}}$            **Participant**
                Forward $S_{ID}$

$\mathbf{S_{ID} > R_{ID}}$            **Non-participant**
                Forward $S_{ID}$    $\rightarrow$ Participant

$\mathbf{S_{ID} < R_{ID}}$            **Non-participant**
                Forward $R_{ID}$    $\rightarrow$ Participant

S$_{ID}$  ID in message from sender

R$_{ID}$  ID at receiver

# Different cases

**S (sender)**                    **R (receiver)**

$$S_{ID} > R_{ID}$$

**Participant**
Forward S$_{ID}$

$$S_{ID} > R_{ID}$$

**Non-participant**
Forward S$_{ID}$        → Participant

$$S_{ID} < R_{ID}$$

**Participant**

$$S_{ID} < R_{ID}$$

**Non-participant**
Forward R$_{ID}$        → Participant

# Different cases

**S (sender)**                          **R (receiver)**

$$S_{ID} > R_{ID}$$

**Participant**
  Forward $S_{ID}$

$$S_{ID} > R_{ID}$$

**Non-participant**
  Forward $S_{ID}$        $\rightarrow$ Participant

$$S_{ID} < R_{ID}$$

**Participant**
  No forwarding (own ID already sent)

$$S_{ID} < R_{ID}$$

**Non-participant**
  Forward $R_{ID}$        $\rightarrow$ Participant

$S_{ID}$  ID in message from sender

$R_{ID}$  ID at receiver

# Different cases

**S (sender)**                    **R (receiver)**

$$S_{ID} > R_{ID}$$

**Participant**
Forward $S_{ID}$

If $S_{ID} = R_{ID}$, it follows

**R** elected as leader

$$S_{ID} > R_{ID}$$

**Non-participant**
Forward $S_{ID}$      $\rightarrow$ Participant

$$S_{ID} < R_{ID}$$

**Participant**
No forwarding (own ID already sent)

$$S_{ID} < R_{ID}$$

**Non-participant**
Forward $R_{ID}$      $\rightarrow$ Participant

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
## Concurrent election start

# Ring-based algorithm:
# Concurrent election start

**Participant**

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
## Concurrent election start



**Participant**

**1**

**Participant**

$1 < 2$, ignore message

**2**

Non-participant

Non-participant

Non-participant

$P_1$

$P_2$

$P_n$

$P_3$

...

$P_j$

# Ring-based algorithm:
# Concurrent election start



**Participant**

**1**

**Participant**

$P_1$

$P_2$

1 < 2, ignore message

**2**

Non-participant

**Participant**

$P_n$

$P_3$

...

$P_j$

Non-participant

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
# Concurrent election start



**Participant**

**1**

**Participant**

P₁  P₂

1 < 2, ignore message

**2**

**Participant**

Pₙ  P₃

**Participant**

**3**

**i+1**

...

Pⱼ  **Participant**

**i**

# Ring-based algorithm:
# Concurrent election start

# Ring-based algorithm:
## Concurrent election start

# Summarizing observations

- **Worst case**: 3N -1 messages
  - *N-1* messages to reach highest ID from lowest ID
  - *N* messages to reach point of origin
  - Leader announcement takes another *N* messages
- **Safety**: even with multiple nodes starting election
- **Liveness**: guaranteed progress *if no failures during the election occur*

# **Self-study questions**

- Develop a mutual exclusion algorithm by reducing the problem to leader election.

- How can the algorithm be made to tolerate 1, 2, … crash faults?

- How can the algorithm be made to tolerate message loss?

## Elections in a Distributed Computing System

HECTOR GARCIA-MOLINA, MEMBER, IEEE

*Abstract*—After a failure occurs in a distributed computing system, it is often necessary to reorganize the active nodes so that they can continue to perform a useful task. The first step in such a reorganization or reconfiguration is to elect a coordinator node to manage the operation. This paper discusses such elections and reorganizations. Two types of reasonable failure environments are studied. For each environment assertions which define the meaning of an election are presented. An election algorithm which satisfies the assertions is presented for each environment.

*Index Terms*—Crash recovery, distributed computing systems, elections, failures, mutual exclusion, reorganization.

I. INTRODUCTION

A DISTRIBUTED system is a collection of autonomous computing nodes which can communicate with each other and which cooperate on a common goal or task [4]. For example, the goal may be to provide the user with a database management system, and in this case the distributed system is called a distributed database [16].

When a node fails or when the communication subsystem which allows nodes to communicate fails, it is usually necessary for the nodes to adapt to the new conditions so that they may continue working on their joint goal. For example, consider a collection of nodes which are processing sensory data and trying to locate a moving object [18]. Each node has some sensors which provide it with a local view of the world. The nodes exchange data and together decide where the object is located. If one of the nodes ceases to operate, the remaining nodes should recognize this and modify their strategy for locating the object. A node which neighbors the failed node could try to collect sensory data for the area which was assigned to the failed node. Another alternative would be for the remaining nodes to use a detection algorithm which is not very sensitive to "holes" in the area being studied. Or the nodes could decide to switch to such an algorithm when the failure occurs. If enough nodes fail, the remaining nodes may decide that they just cannot perform the assigned task, and may select a new or better suited job for themselves.

There are at least two basic strategies by which a distributed system can adapt to failures. One strategy is to have software which can operate continuously and correctly as failures occur and are repaired [9]. (In the previous example, this would correspond to using an algorithm which can detect the object even when there are holes in the data.) The second alternative is to temporarily halt normal operation and to take some time

out to *reorganize* the system. During the reorganization period, the status of the system components can be evaluated, any pending work can either be finished or discarded, and new algorithms (and possibly a new task) that are tailored to the current situation can 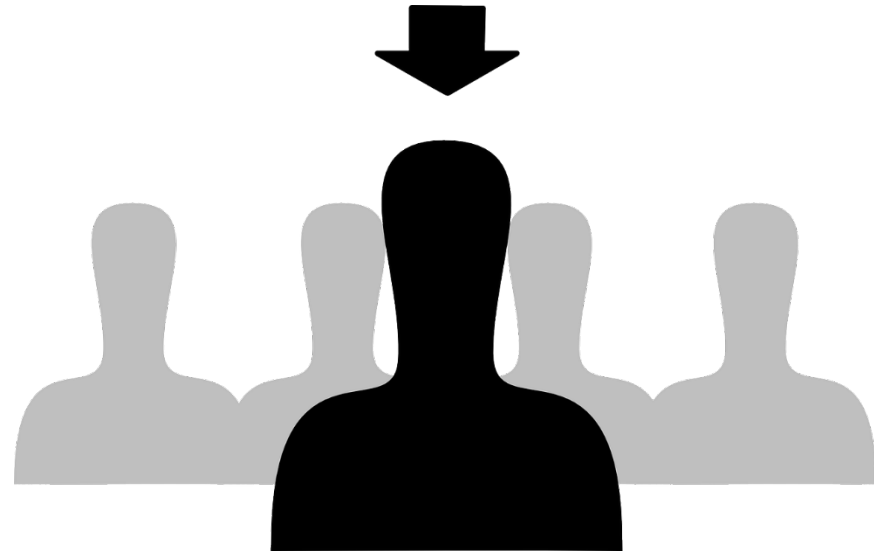be selected. The reorganization of the system is managed by a *single* node called the *coordinator*. (Having more than one node attempting to reorganize will lead to serious confusion.) So as a first step in any reorganization, the operating or active nodes must *elect* a coordinator. It is precisely these elections we wish to study in this paper.

In this paper we will not study the first strategy of continuous operation. This does not mean that we think that the second strategy of reorganization is superior. Which strategy is best depends on the requirements of the application and on the failure rate. If failures are very common, it will probably not pay off to reorganize. If it is not possible to stop performing the task for even a few seconds (or milliseconds?), then clearly continuous operation is a necessity. On the other hand, algorithms for continuous system operation will in all likelihood be more complex than algorithms that can halt when a failure is encountered. Thus, operation between failures should be more efficient if reorganizations are allowed.

In this paper we discuss election protocols in the context of failures, but notice that election protocols can also be used to start up a system initially, or to add or remove nodes from the system [11]. Thus, when the nodes in the system are initially turned on, they will automatically elect a coordinator and start operating, just as if they were recovering from a failure.

The intuitive idea of an election is very natural. We have a number of nodes which talk to each other. After some deliberation among the nodes, a single node is elected the coordinator. When the election terminates, there is only one node that calls itself the coordinator, and all other nodes know the identity of the coordinator. After the election, the coordinator can start the reorganization of the system, after which normal operation can resume.

However, when one attempts to translate the natural idea of an election to a concrete algorithm for performing the election, several interesting issues arise. For example, notice that after a node is elected coordinator, some or all of its constituents may fail. So what does it mean to be coordinator if you really cannot tell who you are coordinating? How can the election protocol cope with failures during the election itself? When certain types of failures occur, it may be impossible to guarantee that a single node calls itself a coordinator. How can these cases be handled? Furthermore, in some situations (like after the communication subsystem is partitioned) we may wish to have more than one coordinator.

Pixabay license

# LEADER ELECTION
## BULLY ALGORITHM, 1982

# Bully algorithm, 1982

- Assumes each node has a **unique ID** (UID), reliable message delivery and synchronous system
- Assumes nodes know each others' UIDs and can directly communicate with one another
  - Higher UIDs have priority
  - Can "bully" nodes with lower UIDs
- Initiated by any node that **detects** leader failure
- Tolerates nodes crashing, even during elections
- Crash recovery model



But I'm the victim here!

# Bully algorithm messages

- ***Election*** message
  - announces an election

- ***Answer*** message
  - responds to an election message

- ***Coordination*** message
  - announces victory and identifies as leader

# @ $P_i$   $P_i$ detects failure of leader

For any $j < i$ and any $i < k$

# @ $P_i$   $P_i$ detects failure of leader



For any *j* < *i* and any *i* < *k*

1. Broadcasts **election message** to all nodes *k* with *k>i*

# @ $P_i$  $P_i$ **detects failure of leader**

For any *j* < *i* and any *i* < *k*

1. Broadcasts **election message** to all nodes *k* with *k>i*
2. Any $P_k$ receiving election message **responds with answer message** and starts another election

# @ $P_i$   $P_i$ **detects failure of leader**

For any *j* < *i* and any *i* < *k*



1. Broadcasts **election message** to all nodes *k* with *k>i*
2. Any $P_k$ receiving election message **responds with answer message** and starts another election
3. Any $P_j$ receiving election message **does not respond**

# @ $P_i$  $P_i$ detects failure of leader

For any *j* < *i* and any *i* < *k*

1. Broadcasts **election message** to all nodes *k* with *k>i*

2. Any $P_k$ receiving election message **responds with answer message** and starts another election

3. Any $P_j$ receiving election message **does not respond**

4. If $P_i$ **does not receive any answer message (timeout)** then it **broadcasts the coordination message**

# @ $P_i$ $P_i$ detects failure of leader

For any **j < i** and any **i < k**

$$j \qquad i \qquad \longrightarrow k$$
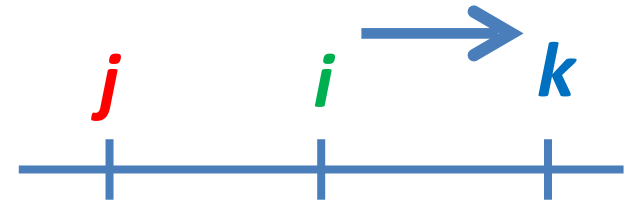
1. Broadcasts **election message** to all nodes *k* with *k>i*
2. Any $P_k$ receiving election message **responds with answer message** and starts another election
3. Any $P_j$ receiving election message **does not respond**
4. If $P_i$ **does not receive any answer message (timeout)** then it **broadcasts the coordination message**
5. If $P_i$ **does receive answer message(s)** then waits to receive **coordination message**

# @ $P_i$   $P_i$ detects failure of leader

For any $j < i$ and any $i < k$

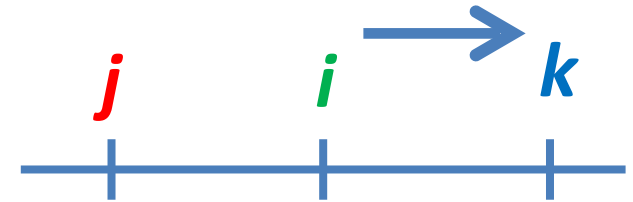$$j \qquad i \qquad k$$

1. Broadcasts **election message** to all nodes $k$ with $k>i$
2. Any $P_k$ receiving election message **responds with answer message** and starts another election
3. Any $P_j$ receiving election message **does not respond**
4. If $P_i$ **does not receive any answer message (timeout)** then it **broadcasts the coordination message**
5. If $P_i$ **does receive answer message(s)** then waits to receive **coordination message**
6. If $P_i$ receives no coordination message before timeout, **it restarts an election**

# @ $P_i$   $P_i$ detects failure of leader

**Special case $i$ = $k$**

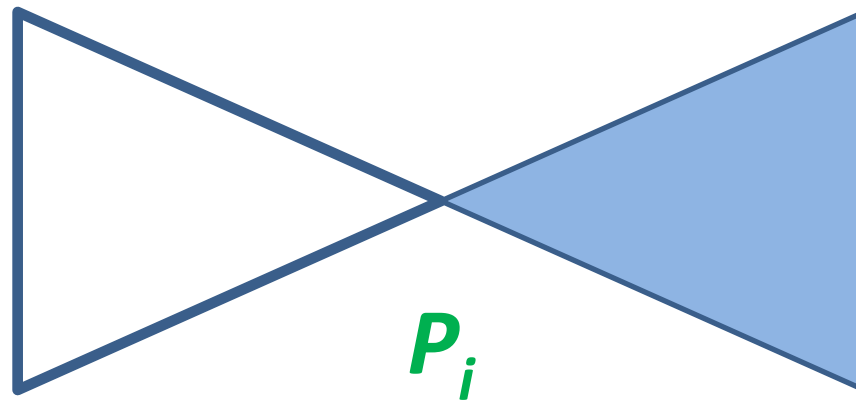For any $j$ < $i$ and any $i$ < $k$

…

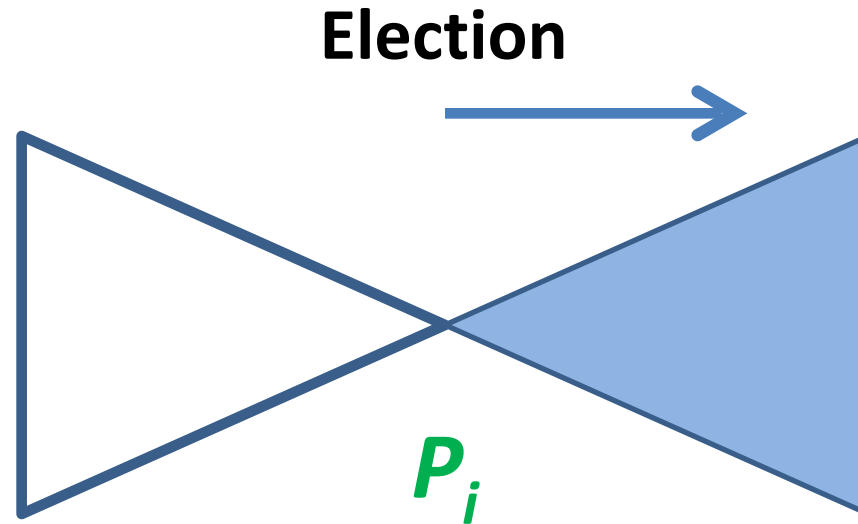4. If $P_i$ **does not receive any answer message (timeout)** then it **broadcasts victory** via **coordination message**

Therefore, if $P_i$ knows it has the highest UID (e.g., after crash recovery), it may as well directly declare itself a leader (sending out the coordinator message).
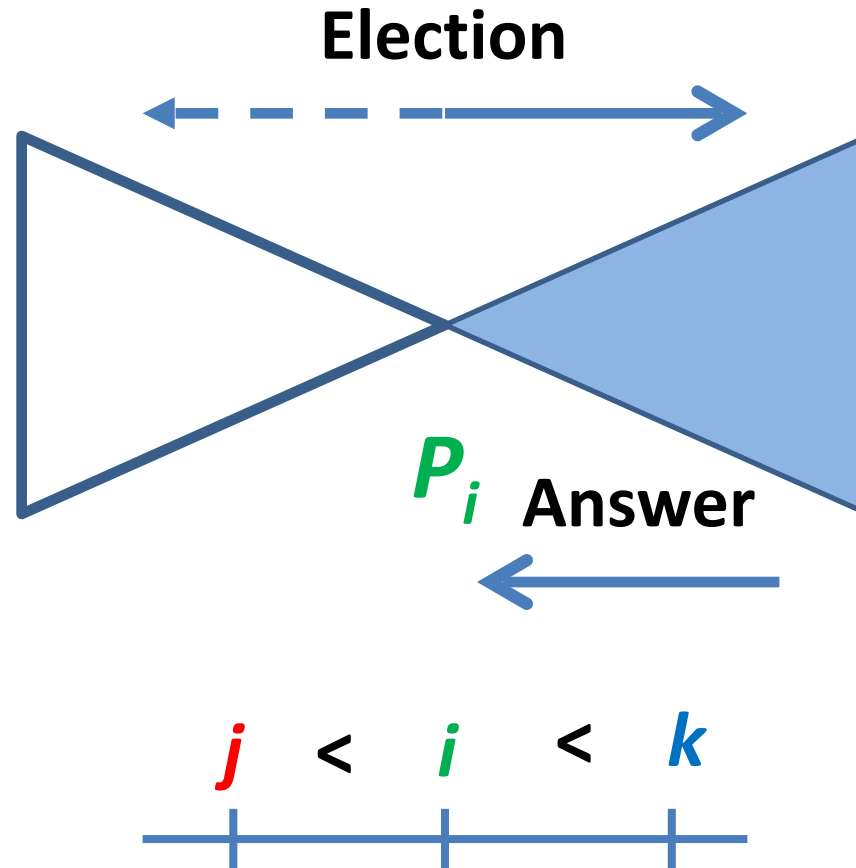
# Election and answer message flow

$$P_i$$

$$j \; < \; i \; < \; k$$

# Election and answer message flow

# Election and answer message flow

# Why Bully?
## Upon node (leader) crash

- Suppose node eventually recovers (*no problem if it stays down, why?*)

- Node may determine that it has the highest UID, thus, it would pronounce itself as leader

  – Even though system may have an existing leader (elected during crash episode)

- New node "**bullies**" current leader

# Summary: Safety and liveness

- **Safety** – argue by contradiction
  - Assume two leaders $P_i$ , $P_k$ (i ≠ k)
  - Thus, $P_i$ , $P_k$ exchanged victory messages
  - Thus, $P_i$ , $P_k$ exchanged election messages
  - But node with lower UID would not have replied
- **Liveness**
  - Would-be leader (*WbL*) fails after sending answer but before sending coordination message
  - Depending on timeout period of nodes with lower UIDs
    - If *WbL* recovers, it issues coordination message
    - If *WbL* does not recover, a new election materializes a leader

# Self-study question

- What are the pros and cons of broadcasting election messages to all nodes vs. to only broadcasting them to nodes with higher UIDs?

- Is the algorithms safe during an election, explain why or why not?

- What is Bully's worst and best case message complexity, assuming *n* nodes in total?