# The Consensus Problem

Roger Wattenhofer

**D**istributed
**C**omputing
**G**roup
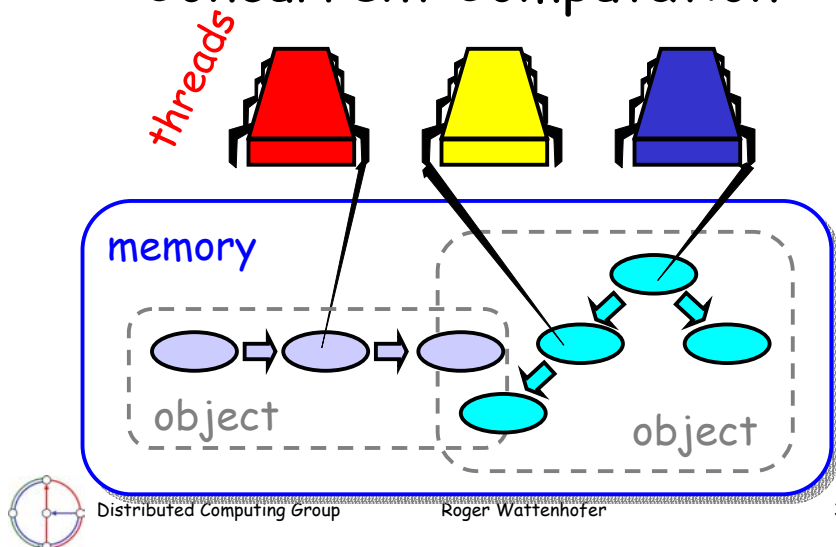
a lot of kudos to
Maurice Herlihy
and Costas Busch
for some of
their slides

---

# Sequential Computation

thread

memory

object

object

---

# Concurrent Computation

threads

memory

object

object

---

# Asynchrony

Sudden unpredictable delays
- Cache misses (*short*)
- Page faults (*long*)
- Scheduling quantum used up (*really long*)

# Model Summary

- Multiple *threads*
  - Sometimes called *processes*
- Single shared *memory*
- *Objects* live in memory
- Unpredictable asynchronous delays

# Road Map

- We are going to focus on principles
  - Start with idealized models
  - Look at a simplistic problem
  - Emphasize correctness over pragmatism
  - "Correctness may be theoretical, but incorrectness has practical impact"

# You may ask yourself …

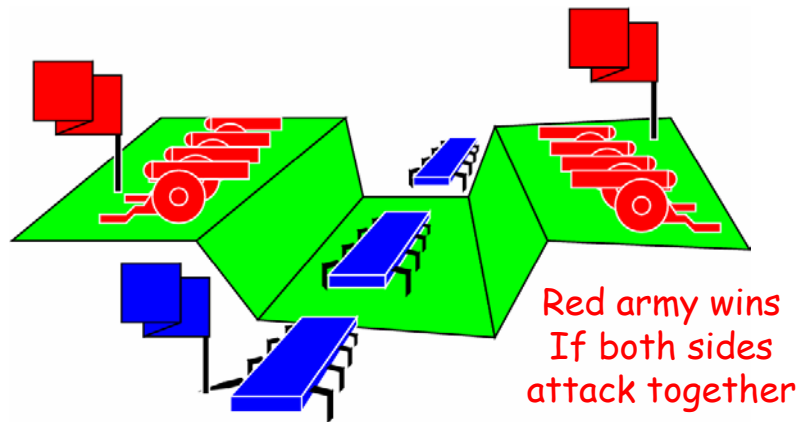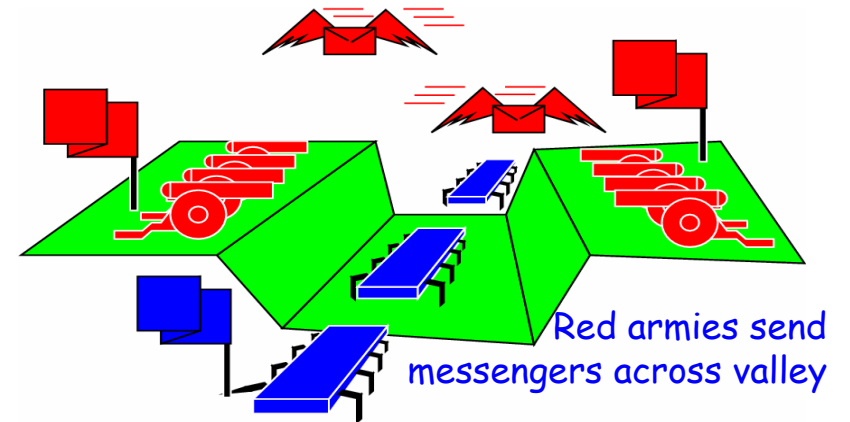I'm no theory weenie - why all the theorems and proofs?

# Fundamentalism

- Distributed & concurrent systems are *hard*
  - Failures
  - Concurrency
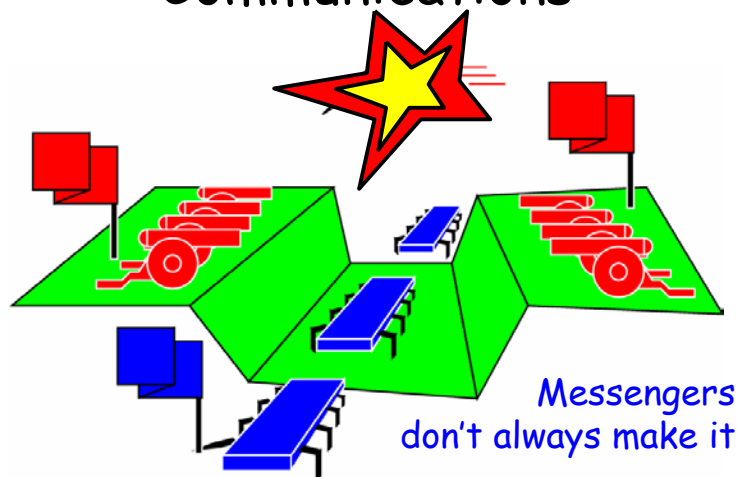- Easier to go from theory to practice than vice-versa

# The Two Generals

Red army wins
If both sides
attack together

# Communications

Red armies send
messengers across valley

# Communications

Messengers
don't always make it

# Your Mission

Design a protocol to ensure
that red armies attack
simultaneously

# Theorem

## There is no non-trivial protocol that ensures the red armies attacks simultaneously

# Proof Strategy

- Assume a protocol exists
- Reason about its properties
- Derive a contradiction

# Proof

1. Consider the protocol that sends fewest messages
2. It still works if last message lost
3. So just don't send it
   - Messengers' union happy
4. But now we have a shorter protocol!
5. Contradicting #1

# Fundamental Limitation

- Need an unbounded number of messages
- Or possible that no attack takes place

## You May Find Yourself ...

I want a real-time YAFA compliant Two Generals protocol using UDP datagrams running on our enterprise-level fiber tachyion network ...

Distributed Computing Group    Roger Wattenhofer    17

## You might say

I want a real-time YAFA ...

p...
r...
fiber tachyion netwo...

Yes, Ma'am, right away!

Distributed Computing Group    Roger Wattenhofer    18

## You might say

Advantage:
• Buys time to find another job
• No one expects software to work anyway

fiber tachyion netwo...

Distributed Computing Group    Roger Wattenhofer    19

## You might say

Advantage:
• Buys time to find another job
• No...
any...

Disadvantage:
• You're doomed
• Without this course, you may not even know you're doomed

Distributed Computing Group    Roger Wattenhofer    20

# You might say

I want a real-time YAFA

I can't find a fault-tolerant algorithm, I guess I'm just a pathetic loser.

fiber tachyon netw

# You might say

**Advantage:**
•No need to take course

I can't find a fault-tolerant algorithm, I guess I'm just a pathetic loser

fiber tachyon netw

# You might say

**Advantage:**
•No need to take course

**Disadvantage:**
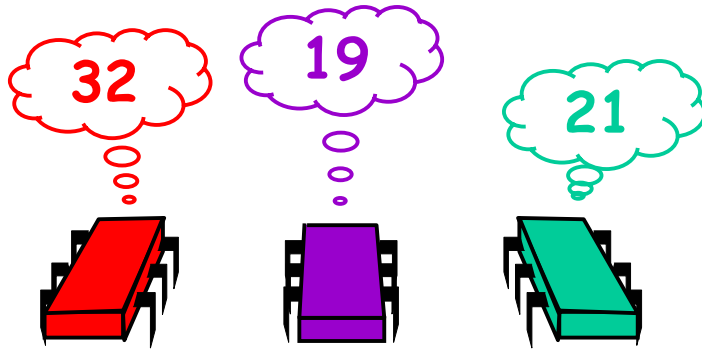•Boss fires you, hires University St. Gallen graduate

# You might say

I want a real-time YAFA

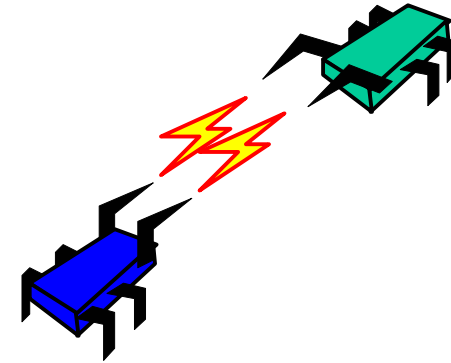Using skills honed in course, I can avert certain disaster!
•Rethink problem spec, or
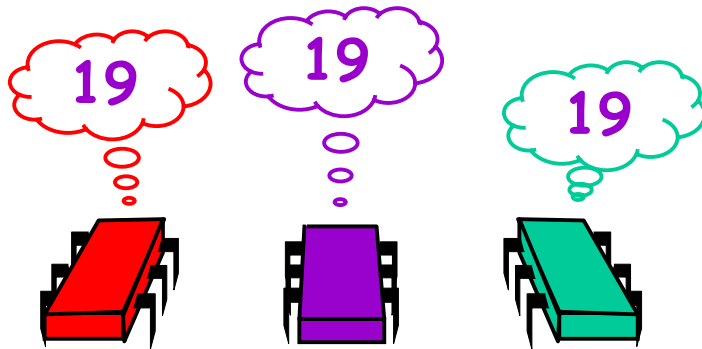•Weaken requirements, or
•Build on different platform

## Consensus: Each Thread has a Private Input

## They Communicate

## They Agree on Some Thread's Input

## Consensus is important

- With consensus, you can implement anything you can imagine…

- Examples: with consensus you can decide on a leader, implement mutual exclusion, or solve the two generals problem

# You gonna learn

- In some models, consensus is possible
- In some other models, it is not

- Goal of this and next lecture: to learn whether for a given model consensus is possible or not … and prove it!

# Consensus #1
## shared memory

- n processors, with n > 1
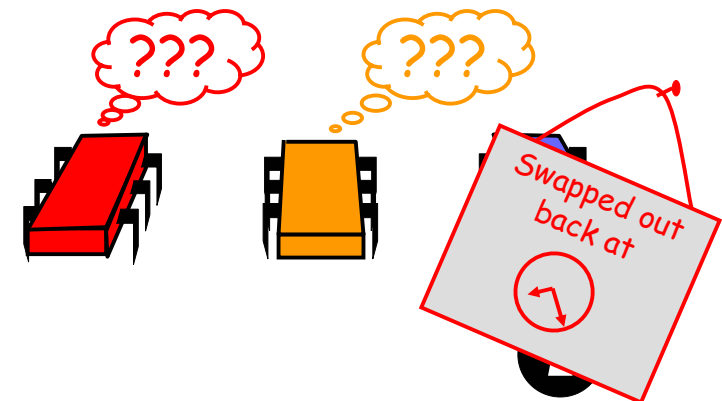- Processors can atomically *read* or *write* (not both) a shared memory cell
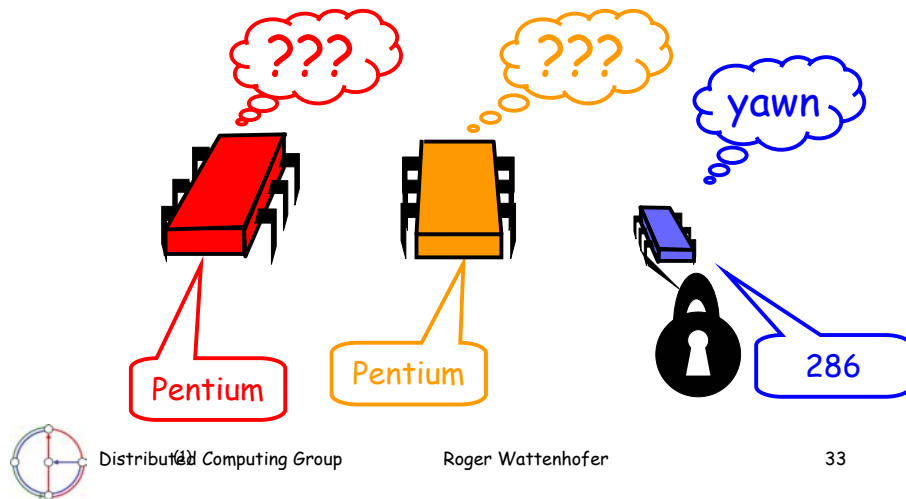
# Protocol (Algorithm?)

- There is a designated memory cell c.
- Initially c is in a special state "?"
- Processor 1 writes its value $v_1$ into c, then decides on $v_1$.
- A processor j (j not 1) reads c until j reads something else than "?", and then decides on that.
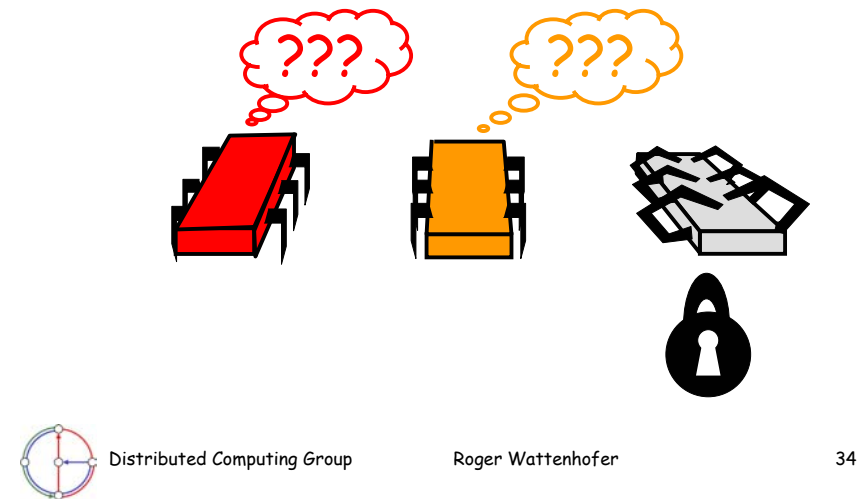
# Unexpected Delay

# Heterogeneous Architectures

# Fault-Tolerance

# Consensus #2
## wait-free shared memory

- n processors, with n > 1
- Processors can atomically *read* or *write* (not both) a shared memory cell
- Processors might crash (halt)
- Wait-free implementation… huh?

# Wait-Free Implementation

- Every process (method call) completes in a finite number of steps
- Implies no mutual exclusion
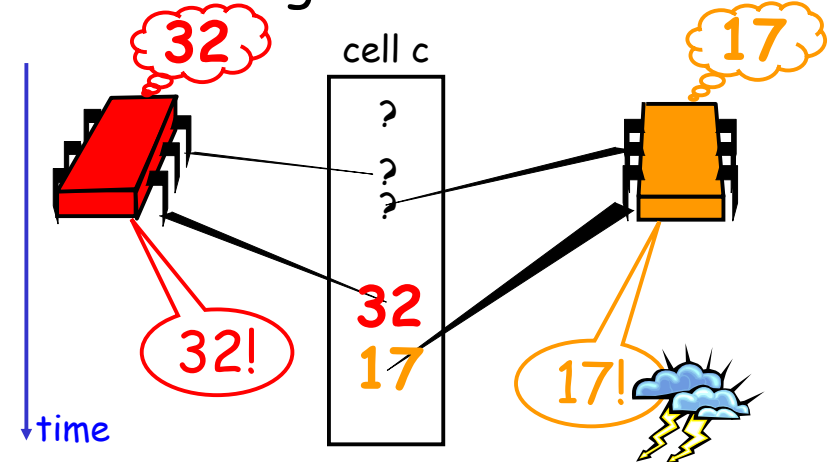- We assume that we have wait-free atomic registers (that is, reads and writes to same register do not overlap)

# A wait-free algorithm…

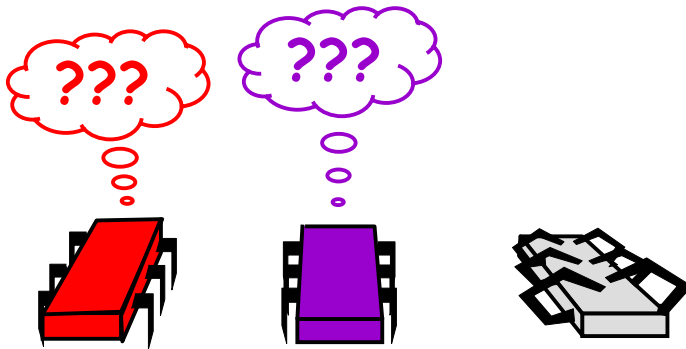- There is a cell c, initially c="?"
- Every processor i does the following

```
r = Read(c);
if (r == "?") then
        Write(c, vᵢ); decide vᵢ;
else
        decide r;
```

# Is the algorithm correct?



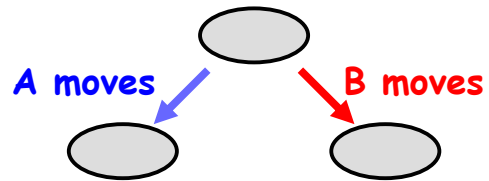cell c

32

17

time

# Theorem: No wait-free consensus

# Proof Strategy

- Make it simple
  - n = 2, binary input
- Assume that there is a protocol
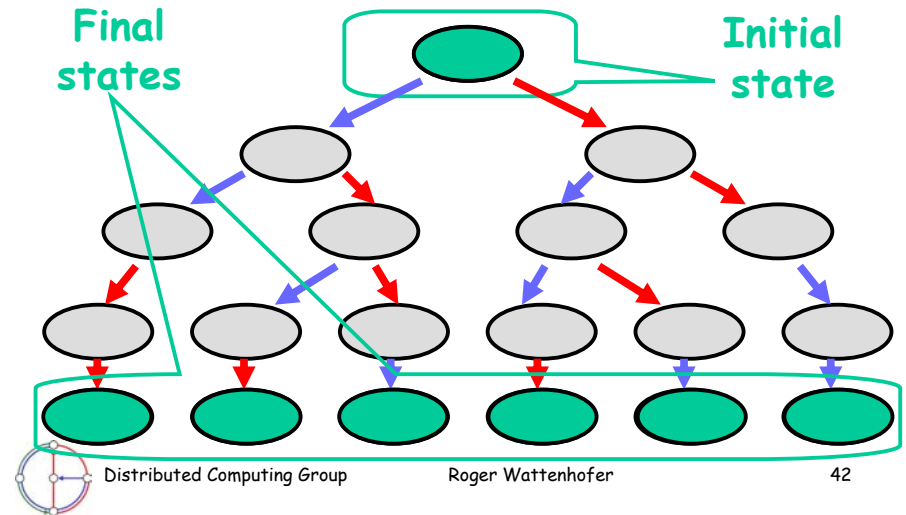- Reason about the properties of any such protocol
- Derive a contradiction
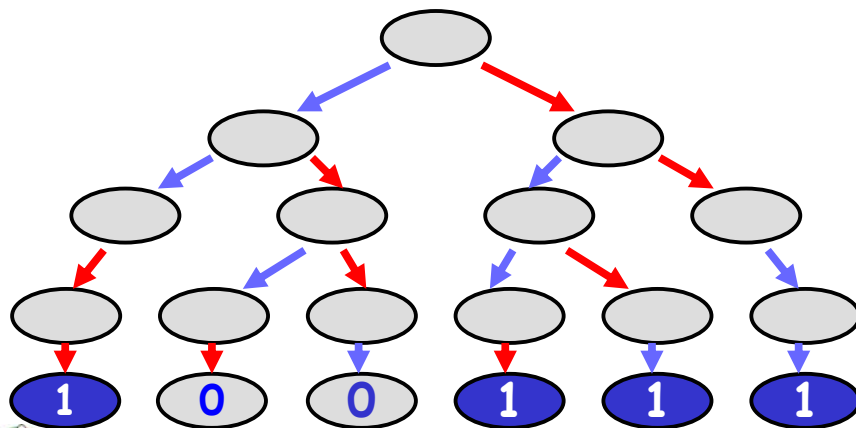
# Wait-Free Computation

**A moves** **B moves**

- Either A or B "moves"
- Moving means
  - Register read
  - Register write

# The Two-Move Tree

**Final states** **Initial state**

# Decision Values

1 0 0 1 1 1

# Bivalent: Both Possible

**bivalent**

1 0 0 1 1 1

# Univalent: Single Value Possible

univalent

1  0  0  1  1  1

# 1-valent: Only 1 Possible

1-valent

1  0  0  1  1  1

# 0-valent: Only 0 possible

0-valent

1  0  0  1  1  1

# Summary

- Wait-free computation is a tree
- Bivalent system states
  - Outcome not fixed
- Univalent states
  - Outcome is fixed
  - May not be "known" yet
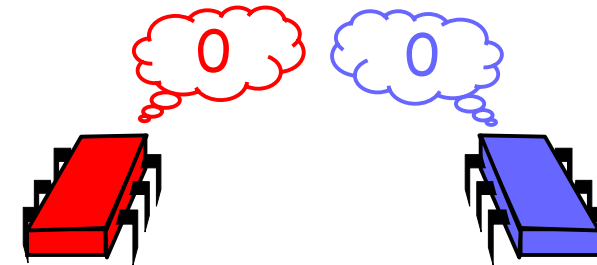  - 1-Valent and 0-Valent states

# Claim

Some initial system state is bivalent

(The outcome is not always fixed from the start.)

# A 0-Valent Initial State

- All executions lead to decision of 0

# A 0-Valent Initial State

- Solo execution by A also decides 0

# A 1-Valent Initial State
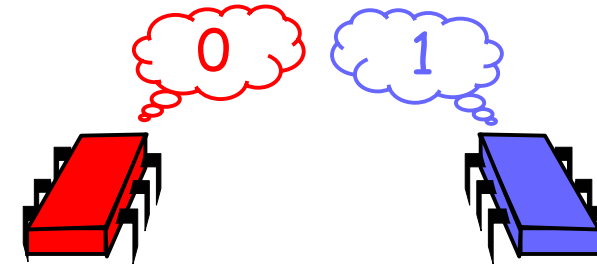
- All executions lead to decision of 1

# A 1-Valent Initial State
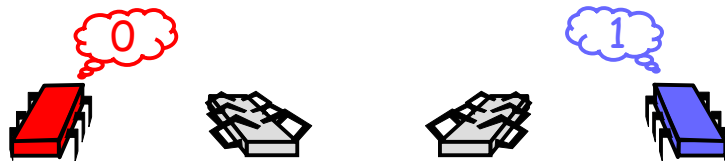


- Solo execution by B also decides 1

# A Univalent Initial State?



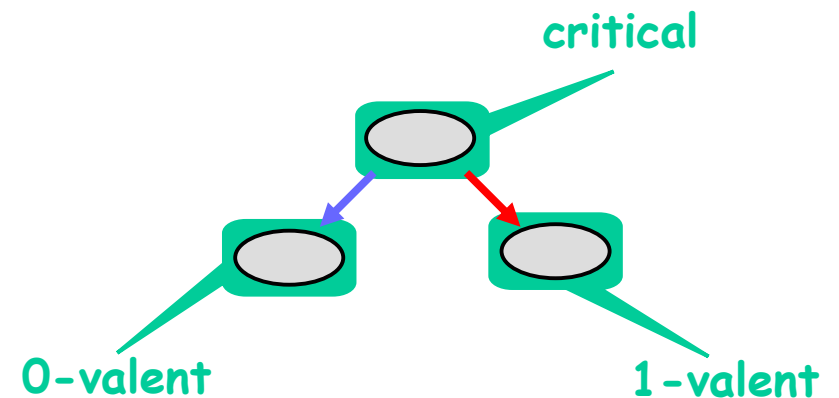- Can all executions lead to the same decision?

# State is Bivalent



- Solo execution by A must decide 0
- Solo execution by B must decide 1

# Critical States

critical



0-valent    1-valent

# Critical States

- Starting from a bivalent initial state
- The protocol can reach a critical state
  - Otherwise we could stay bivalent forever
  - And the protocol is not wait-free

# From a Critical State



**0-valent**          **1-valent**

**If A goes first, protocol decides 0**          **If B goes first, protocol decides 1**

# Model Dependency

- So far, memory-independent!
- True for
  - Registers
  - Message-passing
  - Carrier pigeons
  - Any kind of asynchronous computation

# What are the Threads Doing?

- Reads and/or writes
- To same/different registers

# Possible Interactions

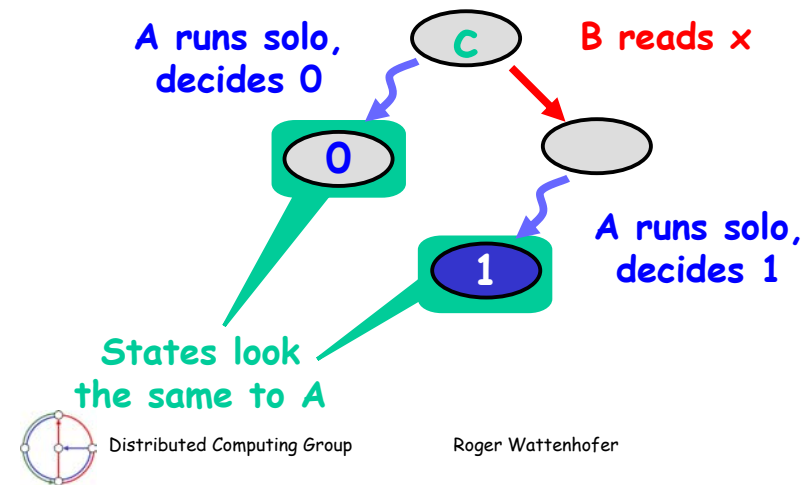|          | x.read() | y.read() | x.write() | y.write() |
|----------|----------|----------|-----------|-----------|
| x.read() | ?        | ?        | ?         | ?         |
| y.read() | ?        | ?        | ?         | ?         |
| x.write()| ?        | ?        | ?         | ?         |
| y.write()| ?        | ?        | ?         | ?         |

# Reading Registers

A runs solo, decides 0

B reads x

C

0

A runs solo, decides 1

1

States look the same to A

# Possible Interactions

|          | x.read() | y.read() | x.write() | y.write() |
|----------|----------|----------|-----------|-----------|
| x.read() | no       | no       | no        | no        |
| y.read() | no       | no       | no        | no        |
| x.write()| no       | no       | ?         | ?         |
| y.write()| no       | no       | ?         | ?         |

# Writing Distinct Registers

A writes y

B writes x

C

B writes x

A writes y

0

1

The song remains the same

## Possible Interactions

|  | x.read() | y.read() | x.write() | y.write() |
|---|---|---|---|---|
| x.read() | no | no | no | no |
| y.read() | no | no | no | no |
| x.write() | no | no | ? | no |
| y.write() | no | no | no | ? |

## Writing Same Registers



C

A writes x

B writes x

A runs solo, decides 0

A writes x

0

1

A runs solo, decides 1

States look the same to A

## That's All, Folks!

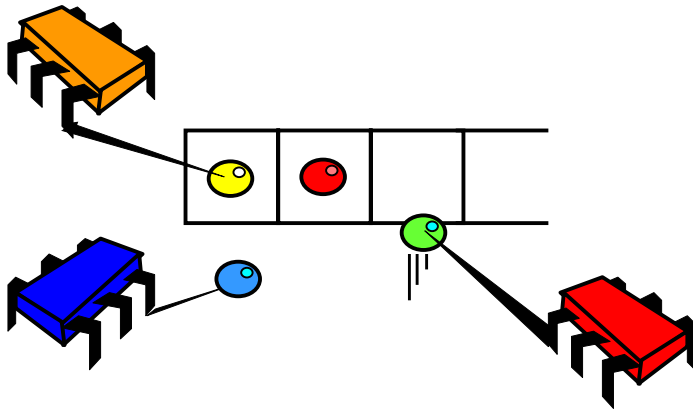|  | x.read() | y.read() | x.write() | y.write() |
|---|---|---|---|---|
| x.read() | no | no | no | no |
| y.read() | no | no | no | no |
| x.write() | no | no | no | no |
| y.write() | no | no | no | no |

## Theorem

- It is impossible to solve consensus using read/write atomic registers
  - Assume protocol exists
  - It has a bivalent initial state
  - Must be able to reach a critical state
  - Case analysis of interactions
    - Reads vs others
    - Writes vs writes

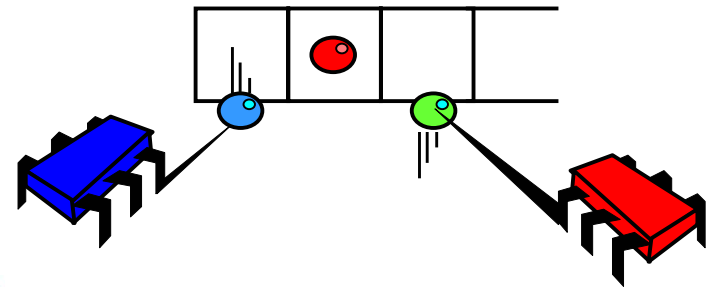# What Does Consensus have to do with Distributed Systems?

# We want to build a Concurrent FIFO Queue

# With Multiple Dequeuers!

# A Consensus Protocol

**2-element array**

**FIFO Queue with red and black balls**

Coveted red ball        Dreaded black ball

# Protocol: Write Value to Array

# Protocol: Take Next Item from Queue

# Protocol: Take Next Item from Queue

**I got the coveted red ball, so I will decide my value**

**I got the dreaded black ball, so I will decide the other's value from the array**

# Why does this Work?

- If one thread gets the red ball
- Then the other gets the black ball
- Winner can take her own value
- Loser can find winner's value in array
  - Because threads write array before dequeuing from queue

# Implication

- We can solve 2-thread consensus using only
  - A two-dequeuer queue
  - Atomic registers

# Implications

- Assume there exists
  - A queue implementation from atomic registers
- Given
  - A consensus protocol from queue and registers
- Substitution yields
  - A wait-free consensus protocol from atomic registers

**contradiction**

# Corollary

- It is impossible to implement a two-dequeuer wait-free FIFO queue with read/write shared memory.

- This was a proof by reduction; important beyond NP-completeness…

# Consensus #3
# read-modify-write shared mem.

- n processors, with n > 1
- Wait-free implementation
- Processors can atomically read *and* write a shared memory cell in one atomic step: the value written can depend on the value read
- We call this a RMW register

# Protocol

- There is a cell c, initially c="?"
- Every processor i does the following

    RMW(c), with

    ```
    if (c == "?") then
        Write(c, vi); decide vi;
    else
        decide c;
    ```

    atomic step

# Discussion

- Protocol works correctly
    - One processor accesses c as the first; this processor will determine decision
- Protocol is wait-free
- RMW is quite a strong primitive
    - Can we achieve the same with a weaker primitive?

# Read-Modify-Write more formally

- Method takes 2 arguments:
    - Variable x
    - Function f
- Method call:
    - Returns value of x
    - Replaces x with f(x)

# Read-Modify-Write

```
public abstract class RMW {
    private int value:

                                    Return prior value
    public void rmw(function f) {
        int prior  = this.value;
        this.value = f(this.value);
        return prior;
    }

}
```

Apply function

# Example: Read

```
public abstract class RMW {
  private int value;

  public void read() {
    int prior  = this.value;
    this.value = this.value;
    return prior;
  }

}
```

**identity function**

# Example: test&set

```
public abstract class RMW {
  private int value;

  public void TAS() {
    int prior  = this.value;
    this.value = 1;
    return prior;
  }

}
```

**constant function**

# Example: fetch&inc

```
public abstract class RMW {
  private int value;

  public void fai() {
    int prior  = this.value;
    this.value = this.value+1;
    return prior;
  }

}
```

**increment function**

# Example: fetch&add

```
public abstract class RMW {
  private int value;

  public void faa(int x) {
    int prior  = this.value;
    this.value = this.value+x;
    return prior;
  }

}
```

**addition function**

# Example: swap

```java
public abstract class RMW {
  private int value;

  public void swap(int x) {
    int prior  = this.value;
    this.value = x;
    return prior;
  }

}
```

**constant function**

# Example: compare&swap

```java
public abstract class RMW {
  private int value;

  public void CAS(int old, int new) {
    int prior  = this.value;
    if (this.value == old)
      this.value = new;
    return prior;
  }

}
```

**complex function**

# "Non-trivial" RMW

- Not simply read
- But
  - test&set, fetch&inc, fetch&add, swap, compare&swap, general RMW
- Definition: A RMW is non-trivial if there exists a value $v$ such that $v \neq f(v)$

# Consensus Numbers (Herlihy)

- An object has consensus number $n$
  - If it can be used
    - Together with atomic read/write registers
  - To implement $n$-thread consensus
    - But not $(n+1)$-thread consensus

# Consensus Numbers

- Theorem
  - Atomic read/write registers have consensus number 1

- Proof
  - Works with 1 process
  - We have shown impossibility with 2

# Consensus Numbers

- Consensus numbers are a useful way of measuring synchronization power

- Theorem
  - If you can implement X from Y
  - And X has consensus number c
  - Then Y has consensus number at least c

# Synchronization Speed Limit

- Conversely
  - If X has consensus number c
  - And Y has consensus number d < c
  - Then there is no way to construct a wait-free implementation of X by Y
- This theorem will be very useful
  - Unforeseen practical implications!

# Theorem

- Any non-trivial RMW object has consensus number at least 2
- Implies no wait-free implementation of RMW registers from read/write registers
- Hardware RMW instructions not just a convenience

# Proof

**Initialized to v**

```
public class RMWConsensusFor2
    implements Consensus {
  private RMW r;

  public Object decide() {
    int i = Thread.myIndex();
    if (r.rmw(f) == v)
      return this.announce[i];
    else
      return this.announce[1-i];
}}
```

**Am I first?**

**Yes, return my input**

**No, return other's input**

---

# Proof

- We have displayed
  - A two-thread consensus protocol
  - Using any non-trivial RMW object

---

# Interfering RMW

- Let F be a set of functions such that for all $f_i$ and $f_j$, either
  - They commute: $f_i(f_j(x))=f_j(f_i(x))$
  - They overwrite: $f_i(f_j(x))=f_i(x)$
- Claim: Any such set of RMW objects has consensus number exactly 2

---

# Examples

- Test-and-Set
  - Overwrite
- Swap
  - Overwrite
- Fetch-and-inc
  - Commute

# Meanwhile Back at the Critical State

A about to apply $f_A$
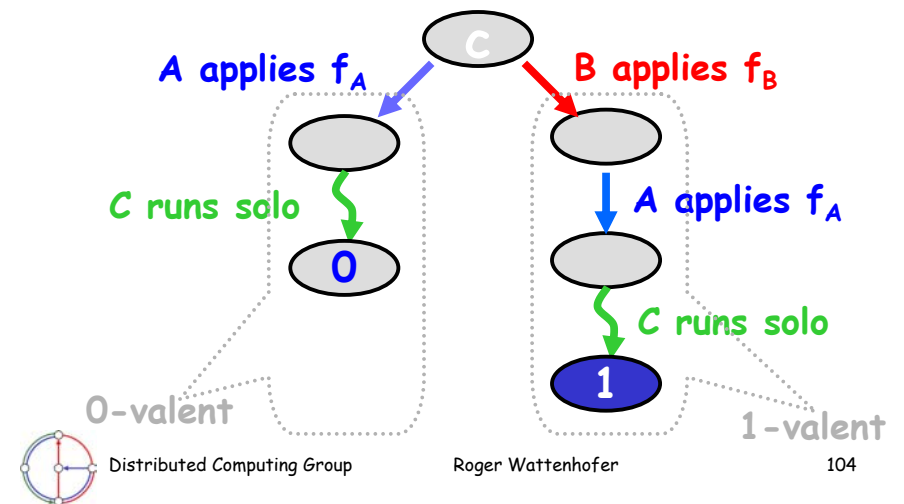
B about to apply $f_B$

C

0-valent

1-valent

# Maybe the Functions Commute

A applies $f_A$

C

B applies $f_B$

B applies $f_B$

A applies $f_A$

C runs solo

C runs solo

0

1

0-valent

1-valent

# Maybe the Functions Commute

A applies $f_A$

C

B applies $f_B$

**These states look the same to C**

B applies $f_B$

A applies $f_A$

C runs solo

C runs solo

0

1

0-valent

1-valent

# Maybe the Functions Overwrite

A applies $f_A$

C

B applies $f_B$

C runs solo

A applies $f_A$

0

C runs solo

1

0-valent

1-valent

# Maybe the Functions Overwrite

**These states look the same to C**

A applies $f_A$ | B applies $f_B$

C runs solo

0

A applies $f_A$

C runs solo

1

0-valent | 1-valent

# Impact

- Many early machines used these "weak" RMW instructions
  - Test-and-set (IBM 360)
  - Fetch-and-add (NYU Ultracomputer)
  - Swap
- We now understand their limitations
  - But why do we want consensus anyway?

# CAS has Unbounded Consensus Number

**Initialized to -1**

**Am I first?**

**Yes, return my input**

**No, return other's input**

```
public class RMWConsensus
    implements Consensus {
  private RMW r;

  public Object decide() {
    int i = Thread.myIndex();
    int j = r.CAS(-1,i);
    if (j == -1)
      return this.announce[i];
    else
      return this.announce[j];
}}
```

# The Consensus Hierarchy

| |
|---|
| **1 Read/Write Registers, …** |
| **2 T&S, F&I, Swap, …** |
| . . . |
| **∞ CAS, …** |

# Consensus #4
## Synchronous Systems

- In real systems, one can sometimes tell if a processor had crashed
  - Timeouts
  - Broken TCP connections

- Can one solve consensus at least in synchronous systems?

# Communication Model
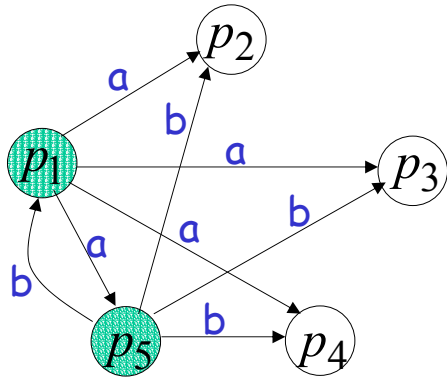
- Complete graph
- Synchronous

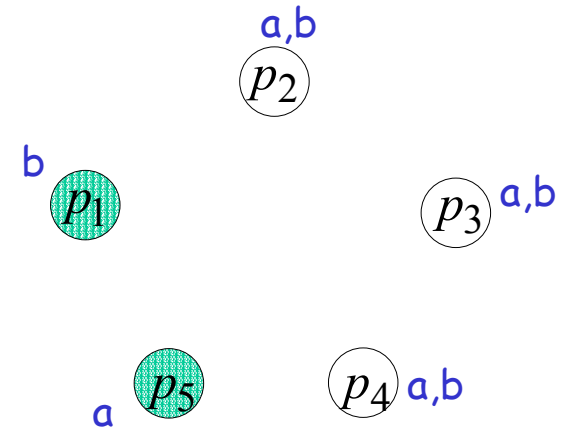# Send a message to all processors in one round: Broadcast
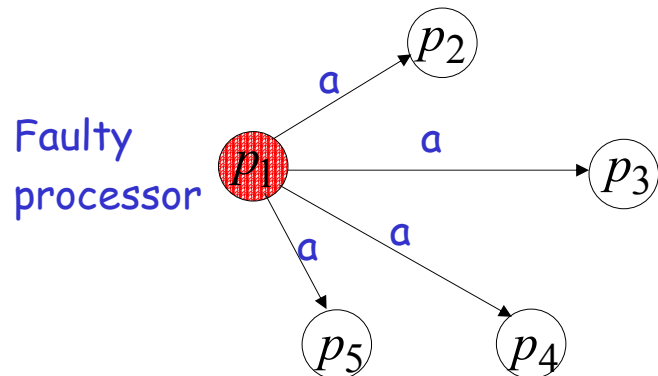
# At the end of the round: everybody receives a

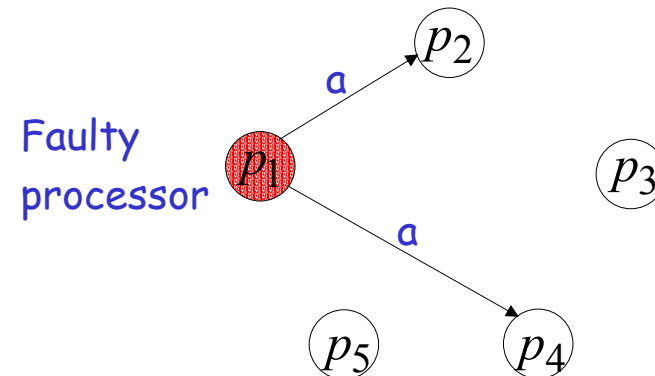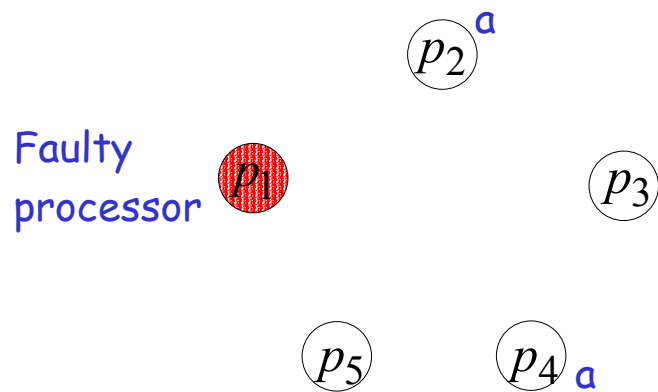## Broadcast: Two or more processes can broadcast in the same round
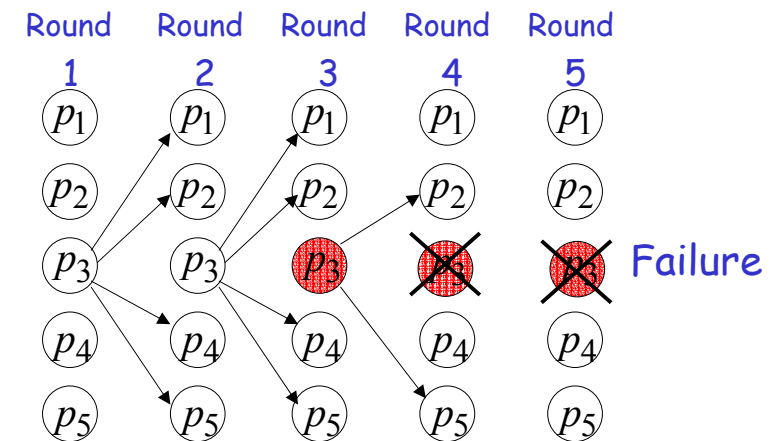
## At end of round...

## Crash Failures

Faulty processor

## Some of the messages are lost, they are never received

Faulty processor

# Effect

Faulty processor

$p_1$

$p_2$ a

$p_3$

$p_5$

$p_4$ a

# After a failure, the process disappears from the network

Round 1   Round 2   Round 3   Round 4   Round 5

$p_1$   $p_1$   $p_1$   $p_1$   $p_1$

$p_2$   $p_2$   $p_2$   $p_2$   $p_2$

$p_3$   $p_3$   $p_3$   $p_3$   $p_3$   Failure

$p_4$   $p_4$   $p_4$   $p_4$   $p_4$

$p_5$   $p_5$   $p_5$   $p_5$   $p_5$

# Consensus:
## Everybody has an initial value

Start

0

1   4

2   3

# Everybody must decide on the same value

Finish

3

3   3

3   3

## Validity condition:

If everybody starts with the same value they must decide on that value
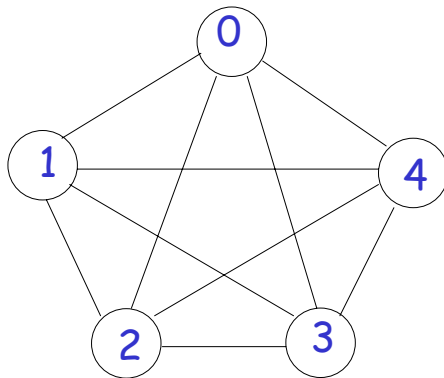
**Start**



**Finish**

---

## A simple algorithm

Each processor:
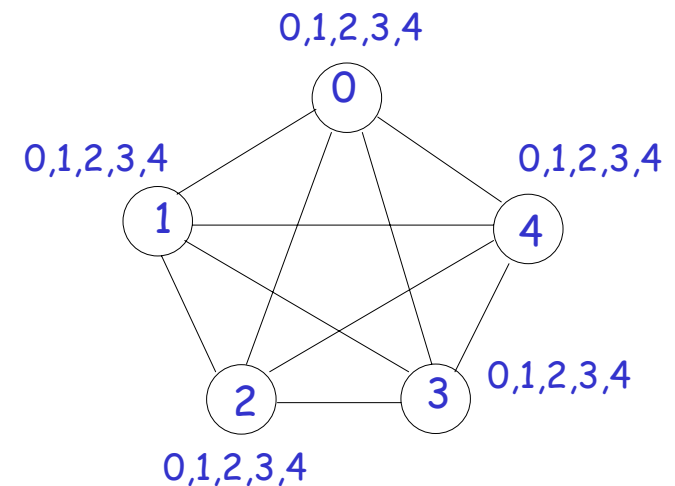
1. Broadcasts value to all processors
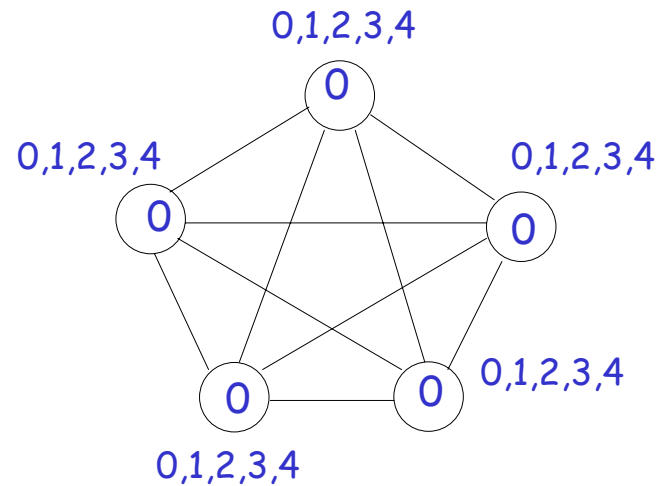
2. Decides on the minimum

(only one round is needed)

---

## Start

---

## Broadcast values

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

# Slide 125

Decide on minimum

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

(graph with all nodes labeled 0)

# Slide 126

Finish

(graph with all nodes labeled 0)

# Slide 127

This algorithm satisfies the validity condition

Start

Finish

(two graphs with all nodes labeled 1)

If everybody starts with the same initial value, everybody sticks to that value (minimum)
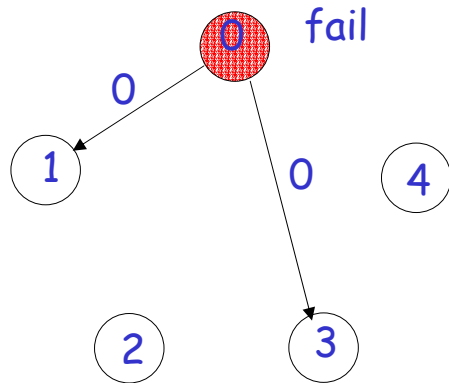
# Slide 128

## Consensus with Crash Failures
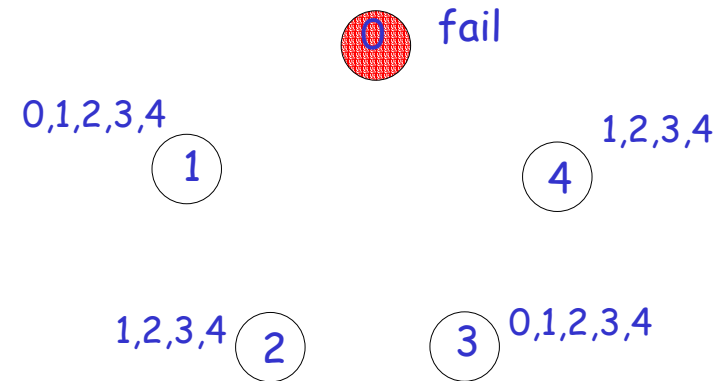
The simple algorithm doesn't work

Each processor:

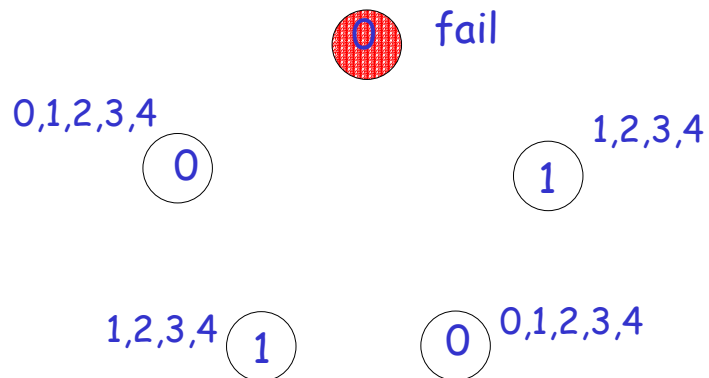1. Broadcasts value to all processors

2. Decides on the minimum

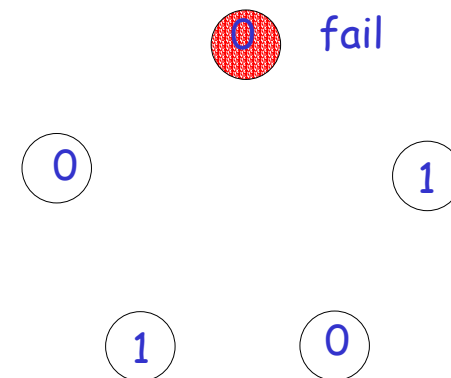## Start — The failed processor doesn't broadcast its value to all processors

0 fail

0 → 1

0 → 3

0

1

4

2

3

## Broadcasted values

0 fail

0,1,2,3,4 — 1

1,2,3,4 — 4

1,2,3,4 — 2

3 — 0,1,2,3,4

## Decide on minimum

0 fail

0,1,2,3,4 — 0

1 — 1,2,3,4

1,2,3,4 — 1

0 — 0,1,2,3,4

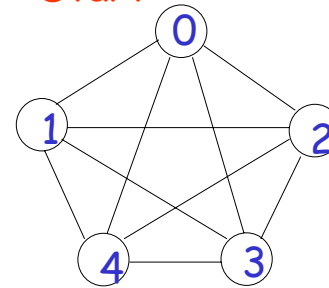## Finish - No Consensus!

0 fail

0

1

1

0

If an algorithm solves consensus for
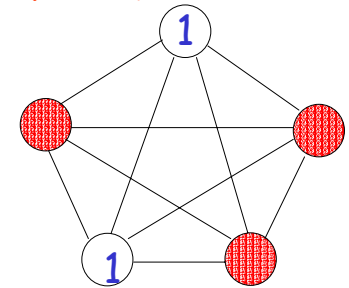f   failed processes we say it is


an f-resilient consensus algorithm

---

Example: The input and output of a
3-resilient consensus algorithm

Start                         Finish
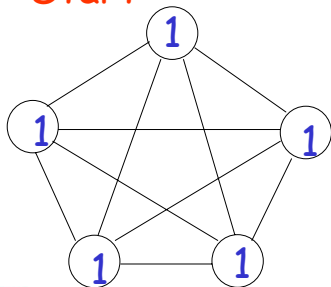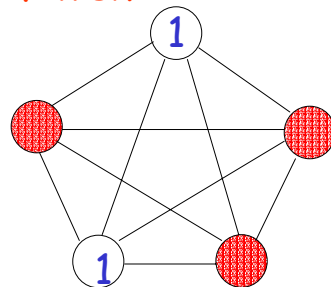
---

New validity condition:

if all non-faulty processes start with the
same value then all non-faulty processes
decide on that value

Start                    Finish

---

An f-resilient algorithm

Round 1:
   Broadcast my value

Round 2 to round f+1:
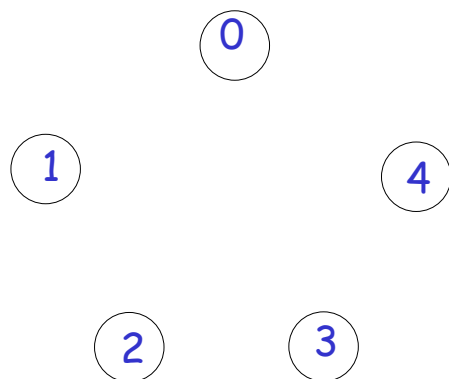   Broadcast any new received values

End of round f+1:
   Decide on the minimum value received
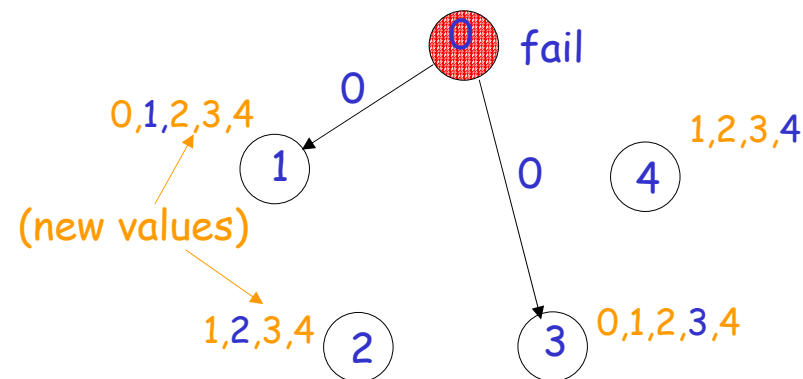
## Slide 137

Example: f=1 failures, f+1=2 rounds needed

Start

## Slide 138

Example: f=1 failures, f+1 = 2 rounds needed

Round 1   Broadcast all values to everybody



0 fail

0,1,2,3,4

1,2,3,4

(new values)

1,2,3,4

0,1,2,3,4

## Slide 139

Example: f=1 failures, f+1 = 2 rounds needed

Round 2   Broadcast all new values to everybody



0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

## Slide 140

Example: f=1 failures, f+1 = 2 rounds needed

Finish   Decide on minimum value



0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

0,1,2,3,4

# Example: f=2 failures, f+1 = 3 rounds needed
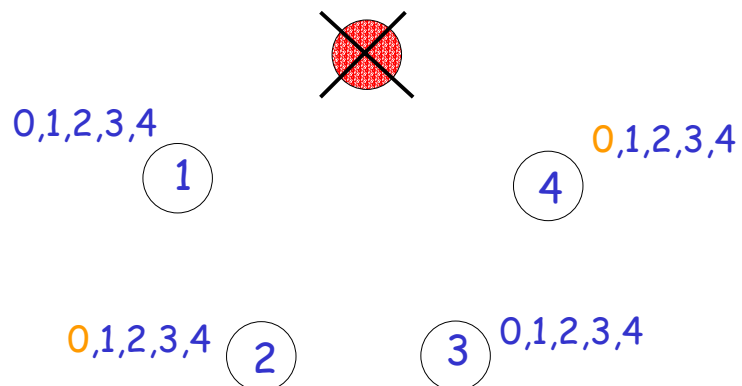
**Start**  Example of execution with 2 failures

# Example: f=2 failures, f+1 = 3 rounds needed

**Round 1**  Broadcast all values to everybody

# Example: f=2 failures, f+1 = 3 rounds needed

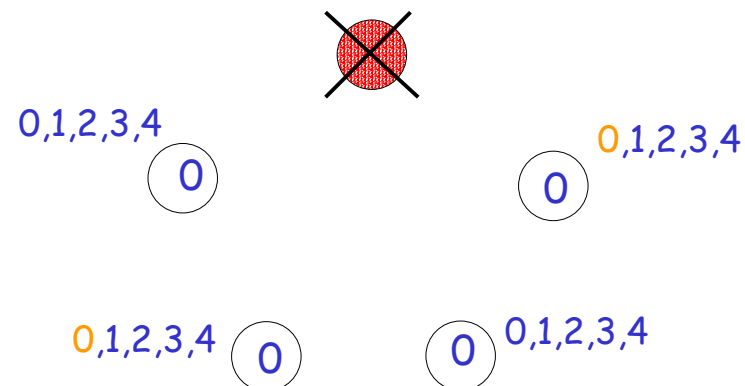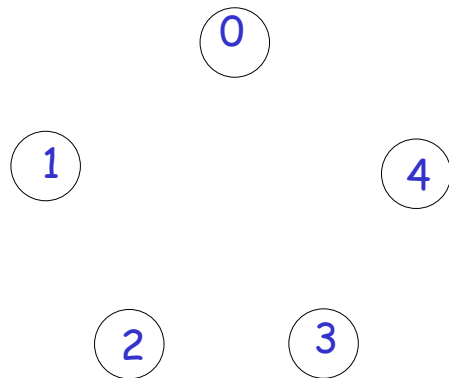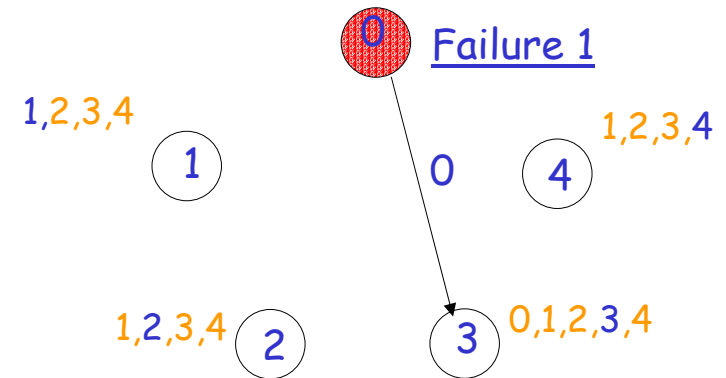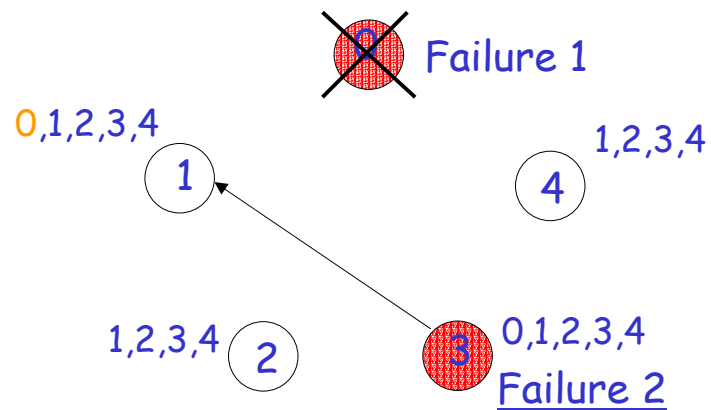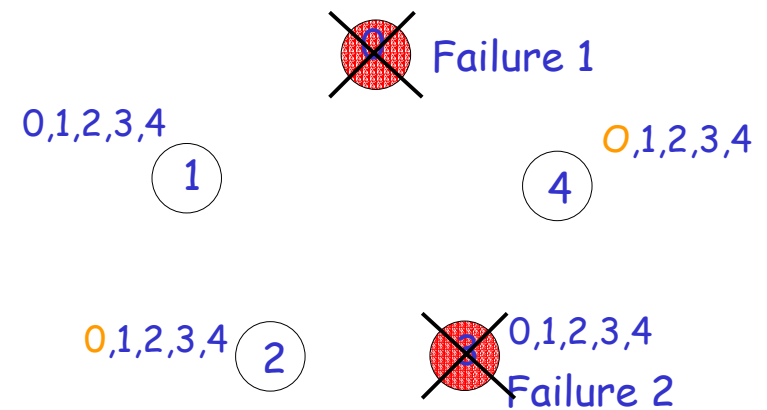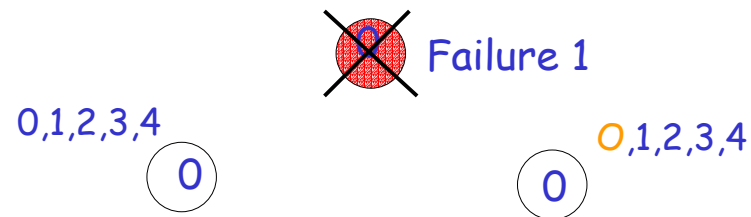**Round 2**  Broadcast new values to everybody

# Example: f=2 failures, f+1 = 3 rounds needed

**Round 3**  Broadcast new values to everybody

**Example: f=2 failures, f+1 = 3 rounds needed**

Finish   Decide on the minimum value



Failure 1

0,1,2,3,4   **0**

O,1,2,3,4   **0**

0,1,2,3,4   **0**

0,1,2,3,4
Failure 2

---

If there are **f** failures and **f+1** rounds then there is a round with no failed process

Round   1   2   3   4   5   6

Example:
5 failures,
6 rounds

No failure

---

# At the end of the round with no failure:

• Every (non faulty) process knows about all the values of all the other participating processes

• This knowledge doesn't change until the end of the algorithm

---

# Therefore, at the end of the round with no failure:

Everybody would decide on the same value

However, as we don't know the exact position of this round,
we have to let the algorithm execute for **f+1** rounds

# Validity of algorithm:

when all processes start with the same
input value then the consensus is that value

This holds, since the value decided from
each process is some input value

# A Lower Bound

Theorem: Any f-resilient consensus algorithm
requires at least f+1 rounds

Proof sketch:

Assume for contradiction that f
or less rounds are enough
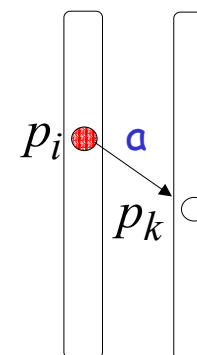
Worst case scenario:

There is a process that fails in
each round

## Worst case scenario

Round     1



before process $p_i$
fails, it sends its
value **a** to only one
process $p_k$

# Worst case scenario

Round  1  2



before process $p_k$
fails, it sends

value **a** to only one
process  $p_m$

---

# Worst case scenario

Round  1  2  3      f



At the end
of round **f**
only one
process $p_n$
knows
about
value **a**

---

# Worst case scenario

Round  1  2  3      f  decide



Process $p_n$
may decide
on **a**, and all
other

processes
may decide
on another
value **(b)**

---

# Worst case scenario

Round  1  2  3      f  decide



Therefore f
rounds are
not enough

At least f+1
rounds are
needed

# Consensus #5
# Byzantine Failures

**Faulty processor**



Different processes receive different values

---

Some messages may be lost

**Faulty processor**



A Byzantine process can behave like a Crashed-failed process

---



Round 1    Round 2    Round 3    Round 4    Round 5    Round 6

Failure                 Failure

After failure the process continues functioning in the network

---

# Consensus with Byzantine Failures

**f-resilient consensus algorithm:**

solves consensus for f failed processes

**Example:** The input and output of a 1-resilient consensus algorithm

Start



Finish

**Validity condition:**

if all non-faulty processes start with the same value then all non-faulty processes decide on that value

Start



Finish

# Lower bound on number of rounds

**Theorem:** Any $f$-resilient consensus algorithm requires at least $f+1$ rounds

**Proof:** follows from the crash failure lower bound

# Upper bound on failed processes

**Theorem:** There is no $f$-resilient algorithm for $n$ processes, where $f \geq n/3$

**Plan:** First we prove the 3 process case, and then the general case

# The 3 processes case

**Lemma:** There is no 1-resilient algorithm for 3 processes

**Proof:** Assume for contradiction that there is a 1-resilient algorithm for 3 processes

---



Local algorithm

Initial value

A(0)
$p_0$

$p_1$    $p_2$

B(1)    C(0)

---



1
$p_0$

$p_1$    $p_2$

1    1

Decision value

---



A(0)    C(1)
$p_3$    $p_2$

B(0)        B(1)
$p_4$        $p_1$

C(0)    A(1)
$p_5$    $p_0$

Assume 6 processes are in a ring

(just for fun)

**Slide 169**

A(0) $p_3$   C(1) $p_2$
B(0) $p_4$   B(1) $p_1$
C(0) $p_5$   A(1) $p_0$

B(1) $p_1$
A(1) $p_0$   C(1)
C(0)   $p_2$
faulty

Processes think they are in a triangle

---

**Slide 170**

A(0) $p_3$   C(1) $p_2$
B(0) $p_4$   B(1) $p_1$
C(0) $p_5$   A(1) $p_0$

1
$p_1$
1
$p_0$   $p_2$
faulty

(validity condition)

---

**Slide 171**

A(0) $p_3$   C(1) $p_2$
B(0) $p_4$   B(1) $p_1$
C(0) $p_5$   A(1) $p_0$

B(0) $p_1$
A(0)   C(0)
$p_0$   $p_2$
A(1)
faulty

1
$p_0$

---

**Slide 172**

A(0) $p_3$   C(1) $p_2$
B(0) $p_4$   B(1) $p_1$
C(0) $p_5$   A(1) $p_0$

0
$p_1$
0
$p_0$   $p_2$
faulty

1
$p_0$

(validity condition)

A(0)   C(1)
$p_3$  $p_2$
B(0)         B(1)
$p_4$        $p_1$
C(0)   A(1)
$p_5$  $p_0$

0   C(0)   A(1)   1
$p_2$  $p_2$  $p_0$  $p_0$
B(0)   B(1)
$p_1$ faulty

A(0)   C(1)
$p_3$  $p_2$
B(0)         B(1)
$p_4$        $p_1$
C(0)   A(1)
$p_5$  $p_0$

0         1
$p_2$   0   1   $p_0$
$p_2$  $p_0$
$p_1$ faulty

# Impossibility

0   1
$p_2$  $p_0$
$p_1$ faulty

# Conclusion

There is no algorithm that solves
consensus for 3 processes
in which 1 is a byzantine process

## The n processes case

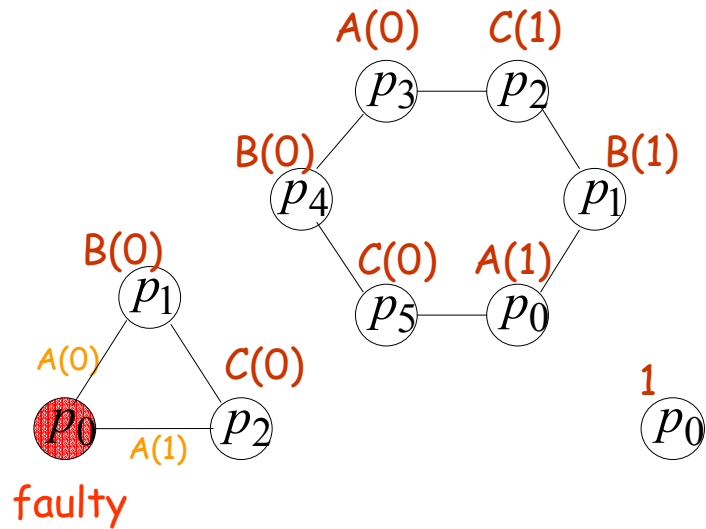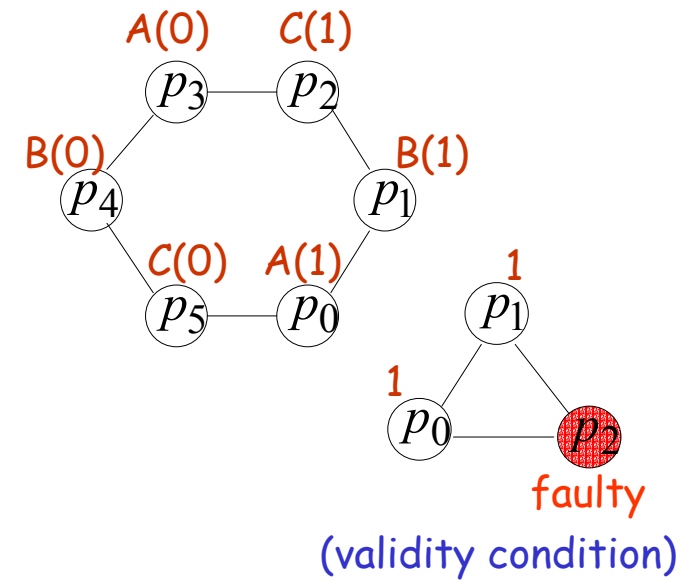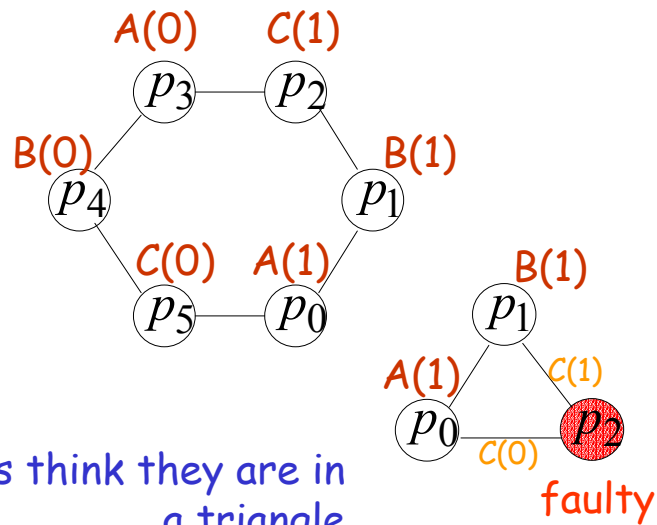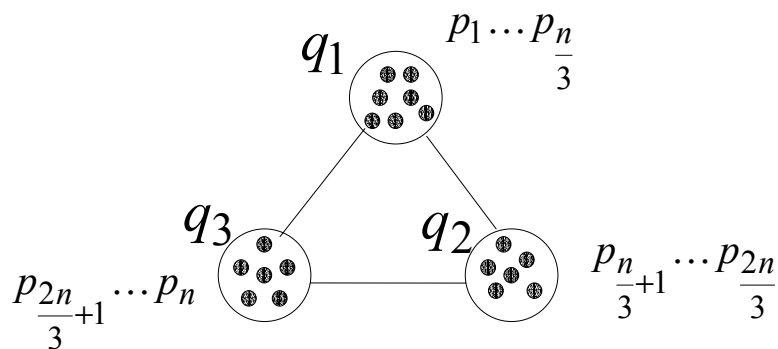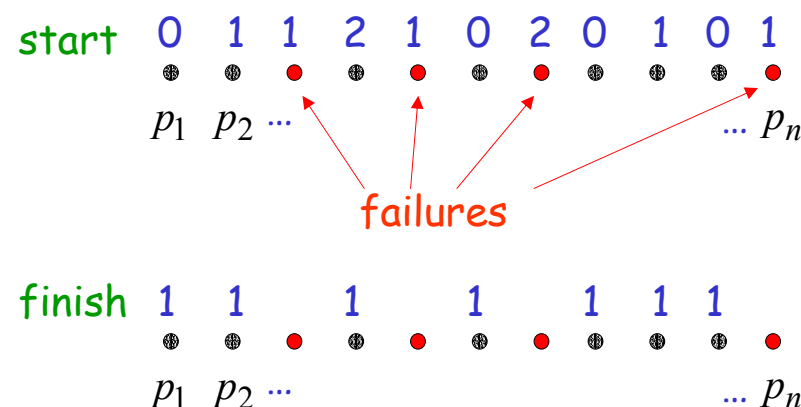Assume for contradiction that there is an $f$-resilient algorithm A for $n$ processes, where $f \geq n/3$
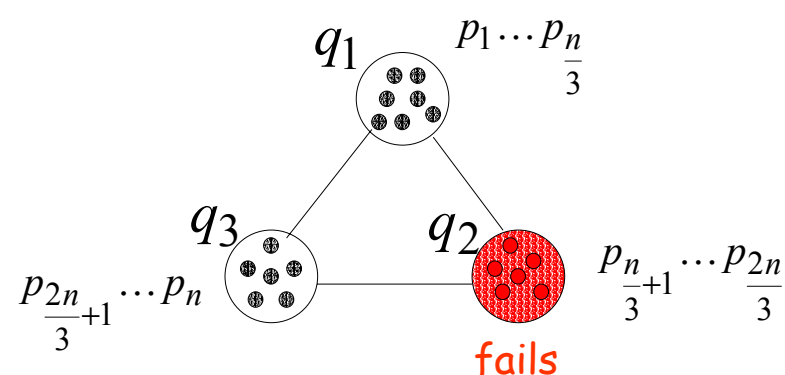
We will use algorithm A to solve consensus for 3 processes and 1 failure (which is impossible, thus we have a contradiction)

## Algorithm A

start  0  1  1  2  1  0  2  0  1  0  1

$p_1$  $p_2$ $\cdots$                                              $\cdots$ $p_n$

failures

finish  1    1         1         1       1   1  1

$p_1$  $p_2$ $\cdots$                                    $\cdots$ $p_n$

$q_1$  $p_1 \ldots p_{\frac{n}{3}}$

$q_3$         $q_2$

$p_{\frac{2n}{3}+1} \ldots p_n$        $p_{\frac{n}{3}+1} \ldots p_{\frac{2n}{3}}$

Each process $q$ simulates algorithm A
on $n/3$ of "$p$" processes

$q_1$  $p_1 \ldots p_{\frac{n}{3}}$

$q_3$         $q_2$

$p_{\frac{2n}{3}+1} \ldots p_n$        $p_{\frac{n}{3}+1} \ldots p_{\frac{2n}{3}}$

fails

When a single $q$ is byzantine, then $n/3$ of
the "$p$" processes are byzantine too.

# Slide 181

Finish of algorithm A

$q_1$

$p_1 \ldots p_{\frac{n}{3}}$

all decide k

$q_3$

$p_{\frac{2n}{3}+1} \ldots p_n$

$q_2$

$p_{\frac{n}{3}+1} \ldots p_{\frac{2n}{3}}$

fails

algorithm A tolerates $n/3$ failures

# Slide 182

Final decision

$q_1$

k

$q_3$

k

$q_2$

fails

We reached consensus with 1 failure

**Impossible!!!**

# Conclusion

There is no $f$-resilient algorithm

for $n$ processes with $f \geq n/3$

# The King Algorithm

solves consensus with $n$ processes and $f$ failures where $f < n/4$ in $f+1$ "phases"
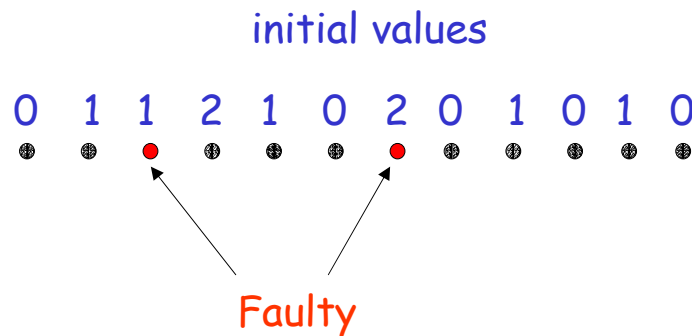
There are $f+1$ phases
Each phase has two rounds
In each phase there is a different king

## Example: 12 processes, 2 faults, 3 kings

initial values

0  1  1  2  1  0  2  0  1  0  1  0

Faulty

## Example: 12 processes, 2 faults, 3 kings

initial values

0  1  1  2  1  0  2  0  1  0  1  0

King 1      King 2      King 3

Remark: There is a king that is not faulty

## The **King** algorithm

Each processor $p_i$ has a preferred value $v_i$

In the beginning, the preferred value
is set to the initial value

## The **King** algorithm: Phase k

Round 1, processor $p_i$ :

- Broadcast preferred value $v_i$

- Set $v_i$ to the majority of
  values received

## The **King** algorithm: Phase k

Round 2, king $p_k$ :

- Broadcast new preferred value $v_k$

Round 2, process $p_i$ :

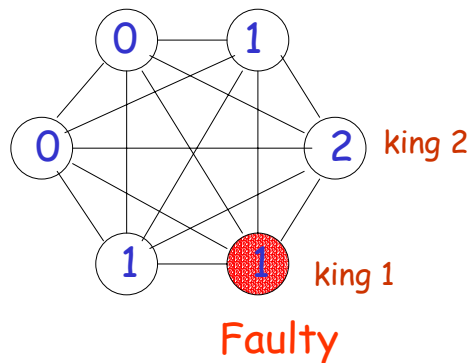- If $v_i$ had majority of less than $\dfrac{n}{2} + f$

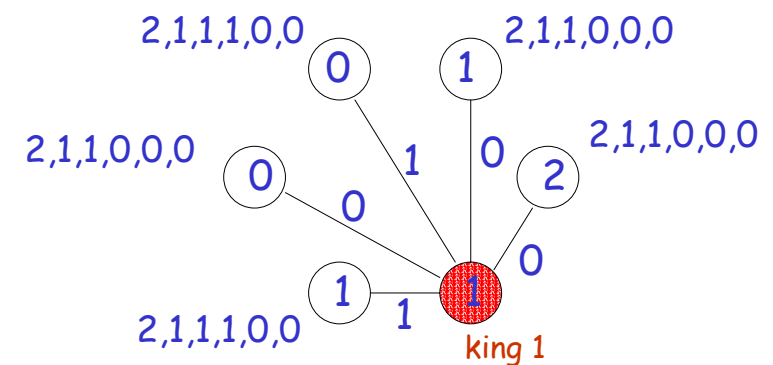  then set $v_i$ to $v_k$

## The **King** algorithm

End of Phase f+1:

Each process decides on preferred value
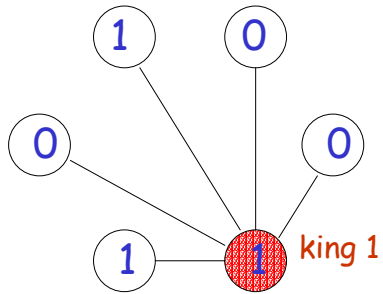
## Example: 6 processes, 1 fault



king 2

king 1

Faulty

## Phase 1, Round 1



2,1,1,1,0,0

2,1,1,0,0,0

2,1,1,0,0,0

2,1,1,0,0,0

2,1,1,1,0,0

king 1

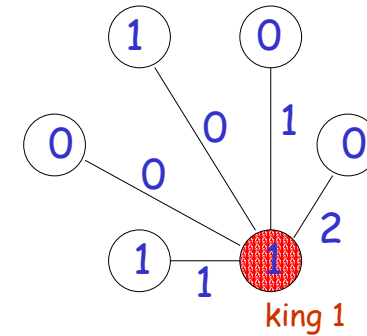Everybody broadcasts

Phase 1, Round 1    Choose the majority

1    0
0         0
1    1 king 1

Each majority population was $3 \le \dfrac{n}{2} + f = 4$
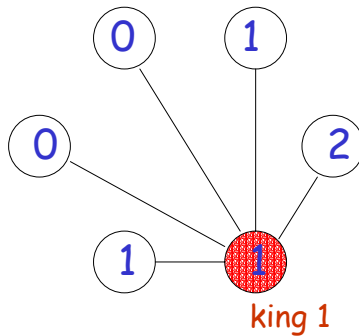
On round 2, everybody will choose the king's value

Phase 1, Round 2

1    0
0    0    1    0
1    1    1    2
king 1

The king broadcasts

Phase 1, Round 2

0    1
0         2
1    1
king 1

Everybody chooses the king's value

Phase 2, Round 1

2,1,1,1,0,0    0    1    2,1,1,0,0,0
2,1,1,0,0,0    0    1    0    2    2,1,1,0,0,0
king 2
1    1    0
2,1,1,1,0,0

Everybody broadcasts

## Phase 2, Round 1    Choose the majority

(1)    (0)

(0)                    (0) king 2

(1)    **(1)** 2,1,1,1,0,0

Each majority population is $3 \leq \dfrac{n}{2} + f = 4$

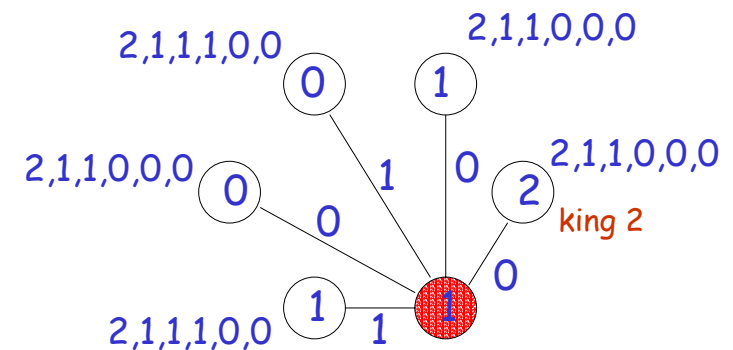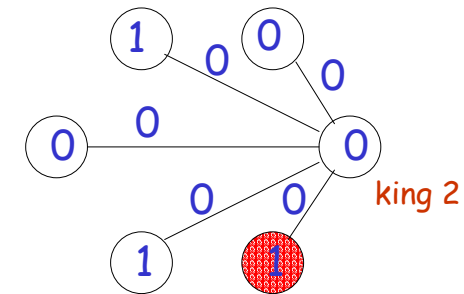On round 2, everybody will choose the king's value

## Phase 2, Round 2

(1)  (0)  (0)  (0)

(0)              (0) king 2

(1)  **(1)**

The king broadcasts

## Phase 2, Round 2

(0)  (0)

(0)        (0) king 2

(0)  **(1)**

Everybody chooses the king's value

Final decision

## Invariant / Conclusion

In the round where the king is non-faulty, everybody will choose the king's value **v**

After that round, the majority will remain value **v** with a majority population which is at least $n - f > \dfrac{n}{2} + f$

# Exponential Algorithm

solves consensus with *n* processes and *f* failures where *f* < *n*/3 in *f* +1 "phases"

But: uses messages with exponential size

# Atomic Broadcast

- One process wants to broadcast message to all other processes
- Either everybody should receive the (same) message, or nobody should receive the message
- Closely related to Consensus: First send the message to all, then agree!

# Summary

- We have solved consensus in a variety of models; particularly we have seen
  - algorithms
  - wrong algorithms
  - lower bounds
  - impossibility results
  - reductions
  - etc.

Questions?

**Distributed**
**Computing**
**Group**
**Roger Wattenhofer**