

# MANY FACES OF CONSENSUS

---

- ▶ What is consensus?
- ▶ Consensus in crash failure model
  - ▶ Scenarios
  - ▶ What is impossible?
  - ▶ What is possible?
- ▶ Consensus in byzantine model
  - ▶ What is impossible?

- ▶ Given a set of parties with their inputs
- ▶ They all must agree on a decision based on their inputs
- ▶ Binary consensus
  - ▶ 0 or 1

Can be extended to any number of decisions

- ▶ Given a set of parties with their inputs
- ▶ They all must agree on a decision based on their inputs
- ▶ Binary consensus
  - ▶ 0 or 1

## CONSISTENCY

All parties  
agree on  
same value

## VALIDITY

Agreed-upon  
value is some  
party's input

## TERMINATION

Each party  
decided in  
finite number  
of steps

# HESITATE AND YOU'RE LOST

---

Crash Failures

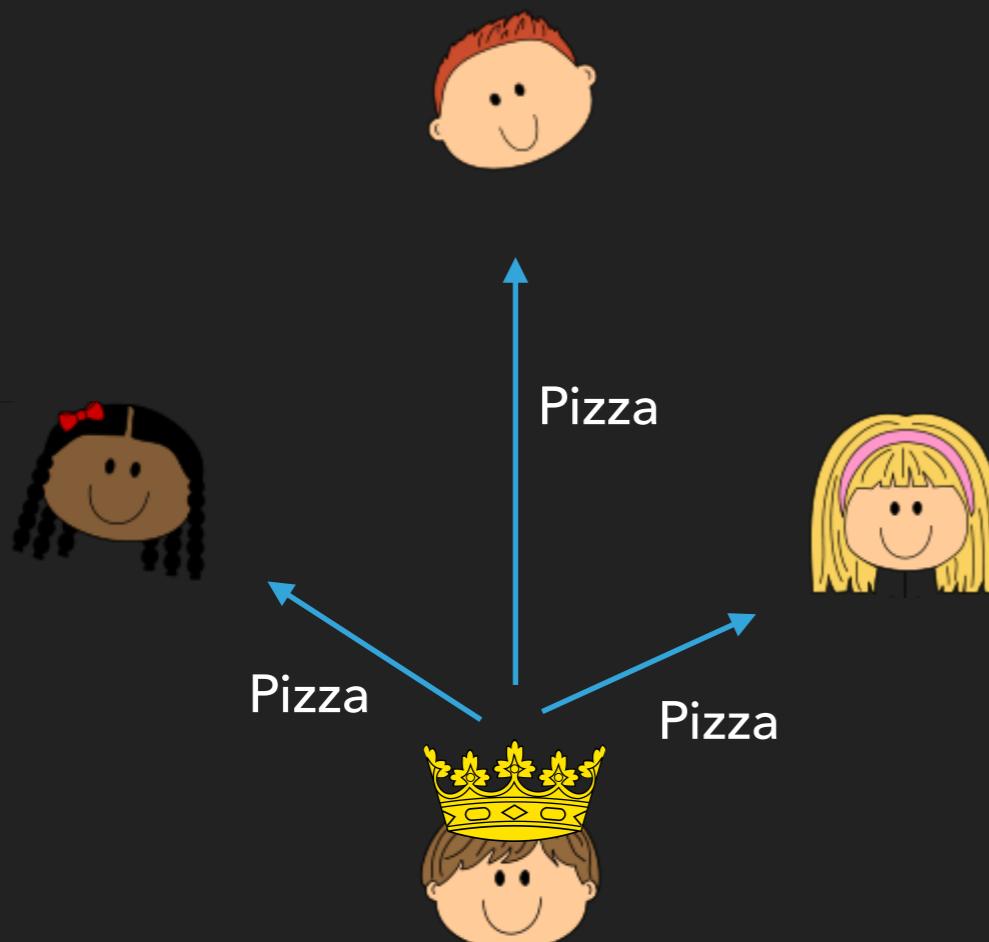
## CASE I



- ▶ Previous communication is allowed
- ▶ Unbounded but reliable communication
- ▶ No one dies

To be or not to be : Pizza or Pasta

## CASE I

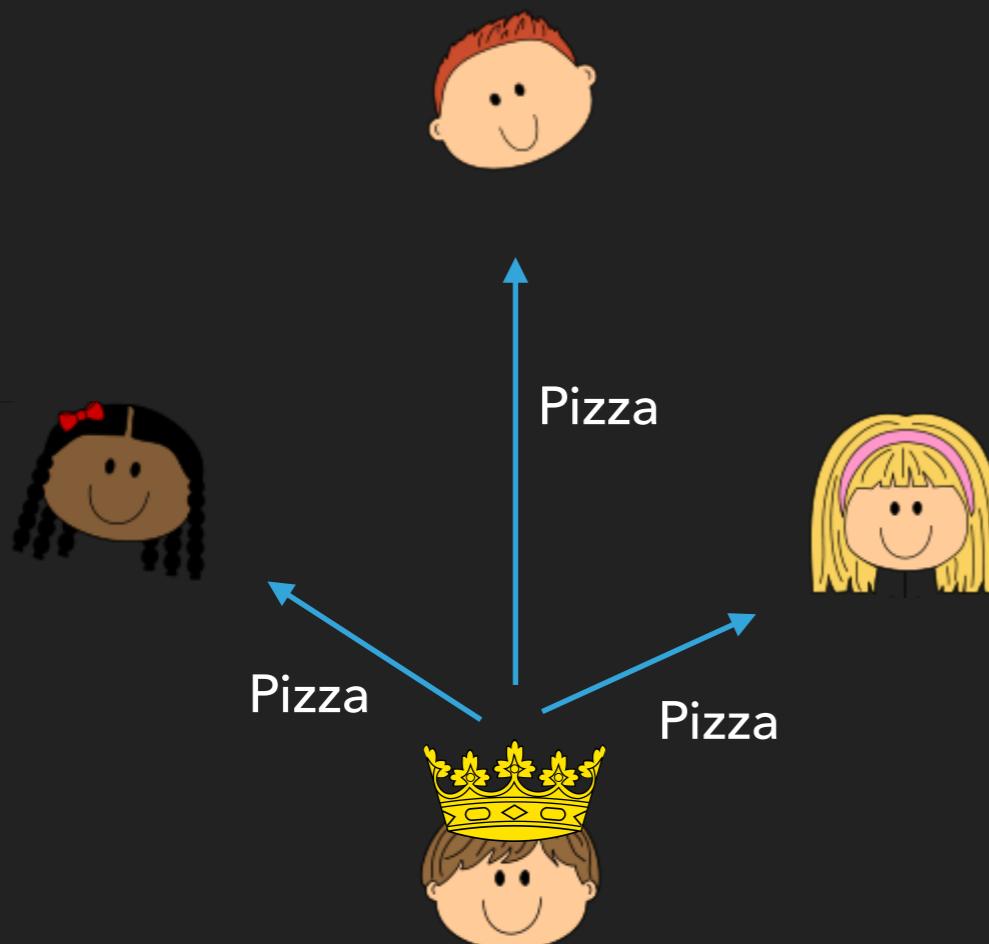


- ▶ Previous communication is allowed
- ▶ Unbounded but reliable communication
- ▶ No one dies

**YES**

To be or not to be : Pizza or Pasta

## CASE II



- ▶ Previous communication is allowed
- ▶ Unbounded but **unreliable** communication
- ▶ No one dies

To be or not to be : Pizza or Pasta

## CASE II



I MUST ACK



DOES HE  
KNOW I RECEIVED  
THE ACK

- ▶ Previous communication is allowed
- ▶ Unbounded but **unreliable** communication
- ▶ No one dies

To be or not to be : Pizza or Pasta

## CASE II



I MUST ACK

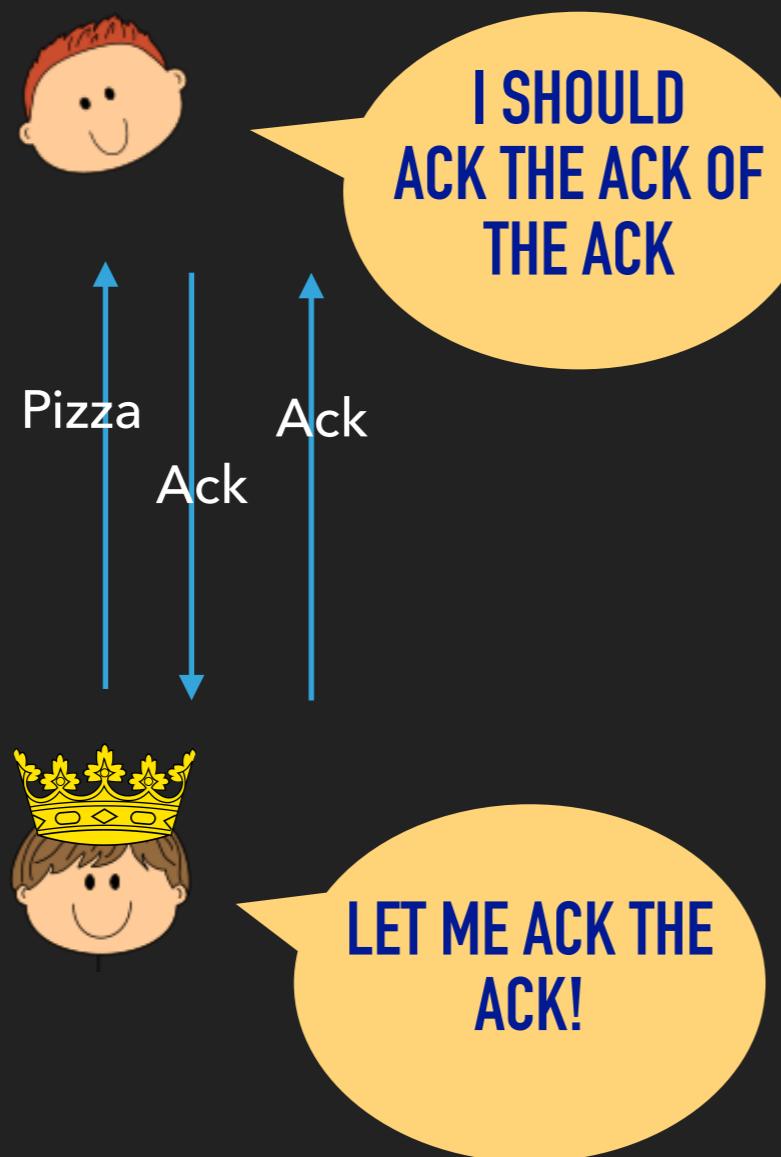


LET ME ACK THE  
ACK!

- ▶ Previous communication is allowed
- ▶ Unbounded but **unreliable** communication
- ▶ No one dies

To be or not to be : Pizza or Pasta

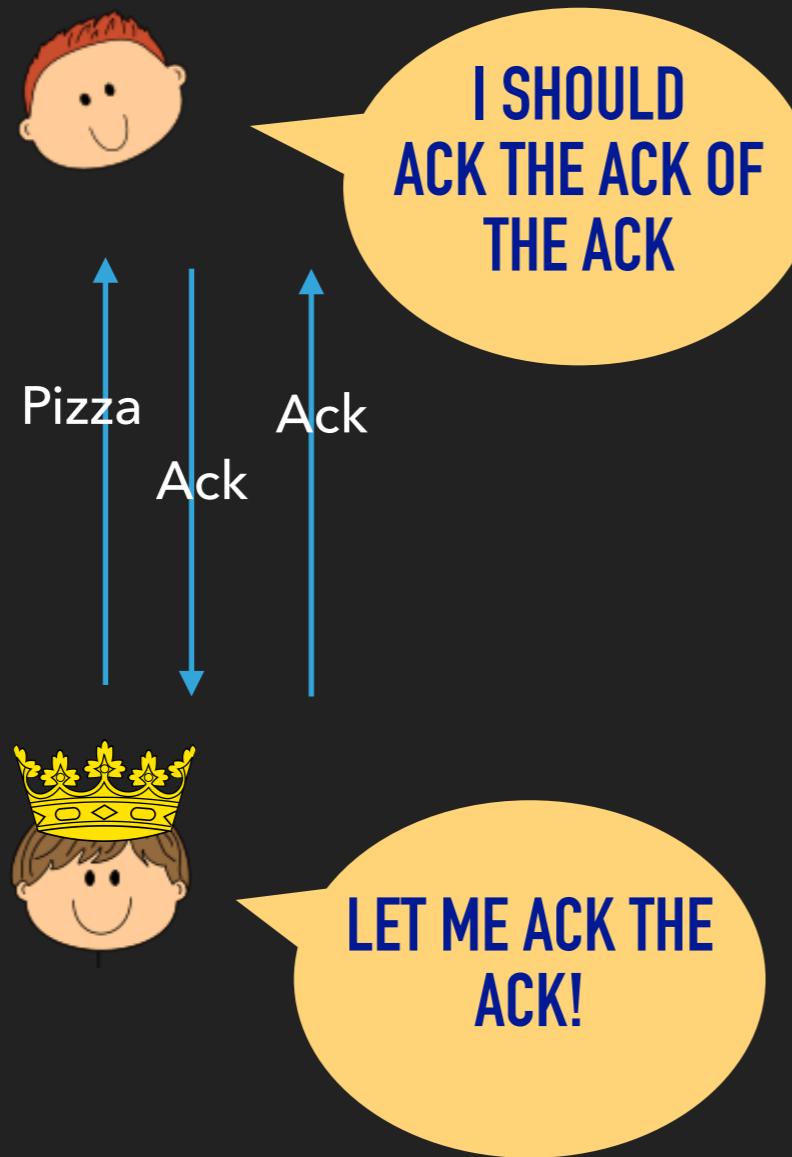
## CASE II



- ▶ Previous communication is allowed
- ▶ Unbounded but **unreliable** communication
- ▶ No one dies

To be or not to be : Pizza or Pasta

## CASE II



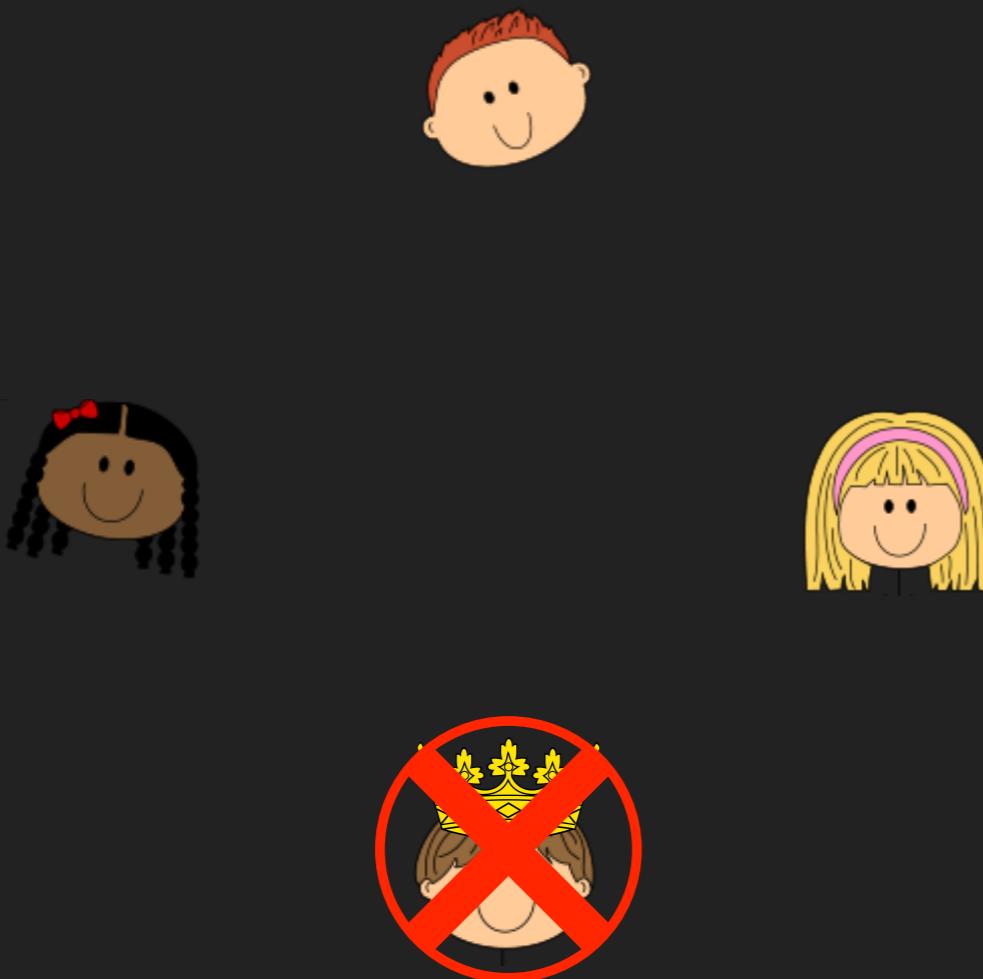
- ▶ Previous communication is allowed
- ▶ Unbounded but **unreliable** communication
- ▶ No one dies

**MAYBE NOT**

QUACK ... QUACK ... QUACK ... QUACK ... QUACK ... QUACK



## CASE III

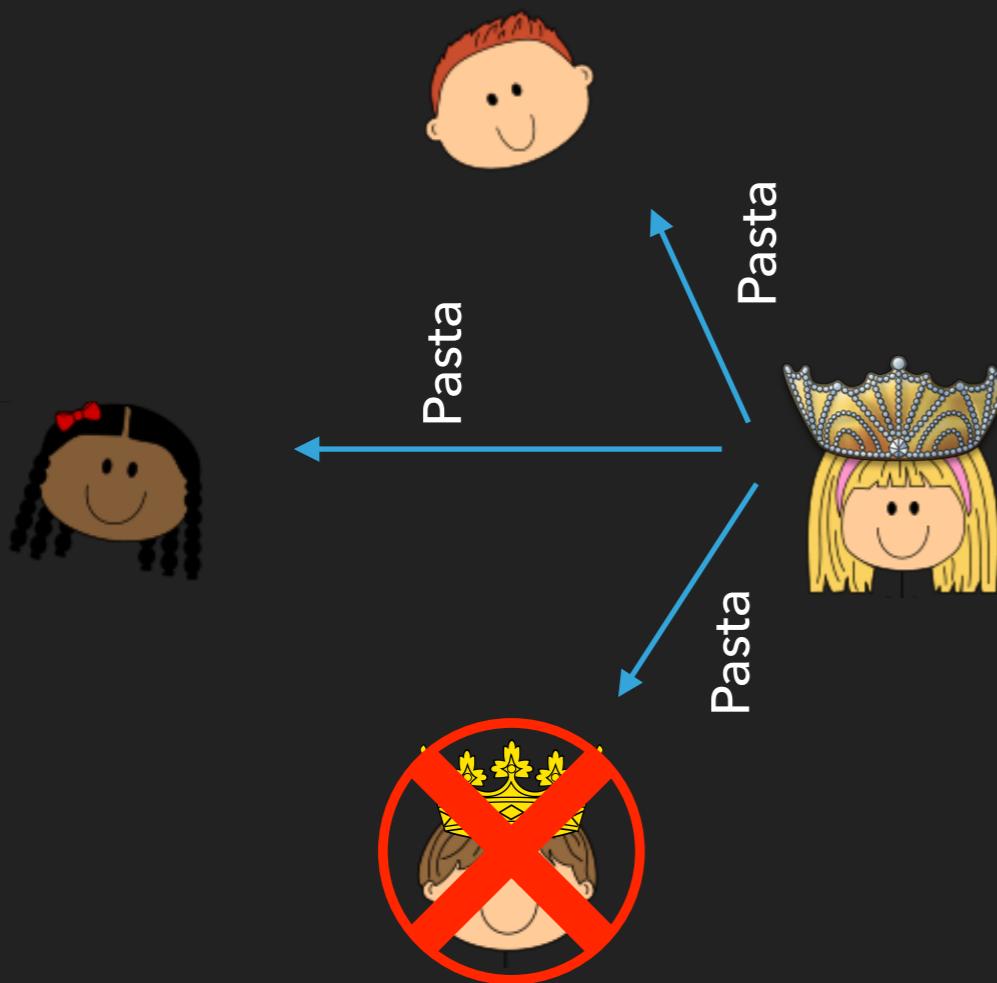


To be or not to be : Pizza or Pasta

- ▶ Previous communication is allowed
- ▶ Unbounded but reliable communication
- ▶ One can die
- ▶ Others can not differentiate b/w whether crown is dead or network is slow

**MAYBE NOT**

## CASE IV



- ▶ Previous communication is allowed
- ▶ **Bounded** but reliable communication
- ▶ One can die

**Use timeouts**

**YES**

To be or not to be : Pizza or Pasta

CONSENSUS DEPENDS  
HEAVILY ON SYSTEM MODEL

Takeaway

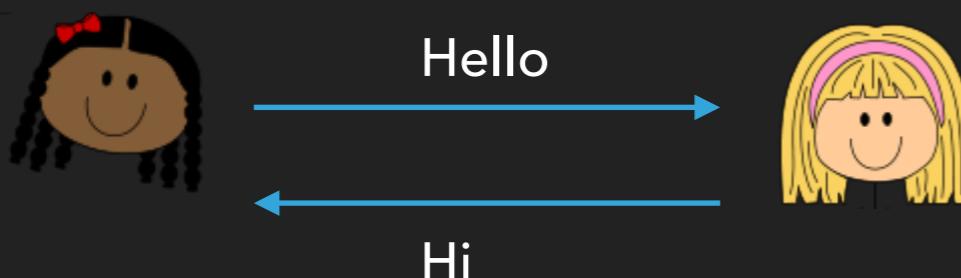
CONSENSUS\* IS IMPOSSIBLE IN A  
COMPLETELY ASYNCHRONOUS  
MESSAGE PASSING SYSTEMS WHERE  
EVEN A SINGLE PROCESSOR FAILS

[Fischer, Lynch, Paterson '85]

# MODELS OF COMMUNICATION

17

## MESSAGE PASSING



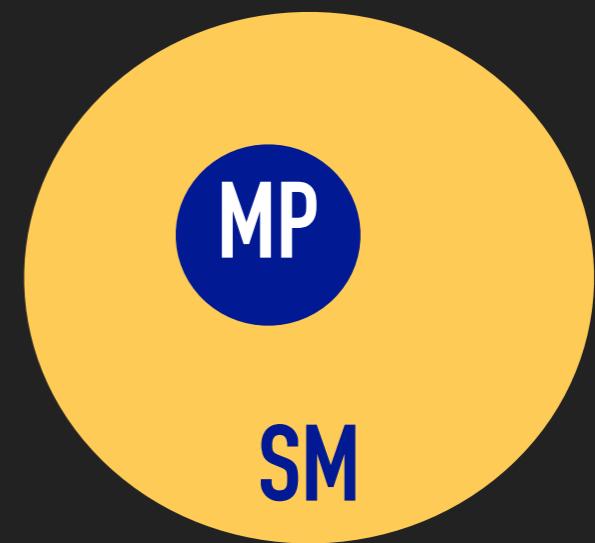
## SHARED MEMORY



FLP

OUR PROOF

- ▶ Proof in asynchronous shared-memory [Herlihy '88]
- ▶ More powerful than message-passing
  - ▶ Inherent broadcasting capabilities
  - ▶ On crashes, value still in memory



**DETOUR**

**MP = SM ( $t < 1/2$ )**

**[Attiya, Noy, Dolev '90]**

Assume otherwise ...

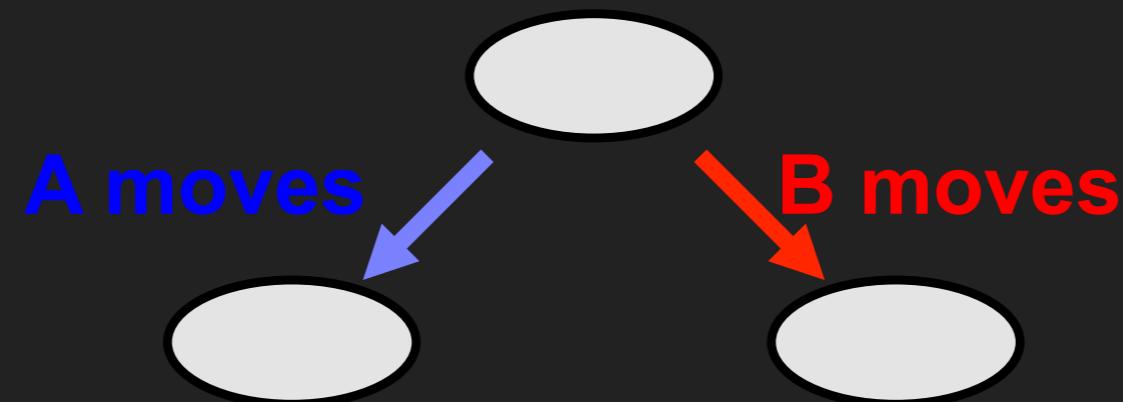
Reason about the properties of any such protocol

Derive a contradiction

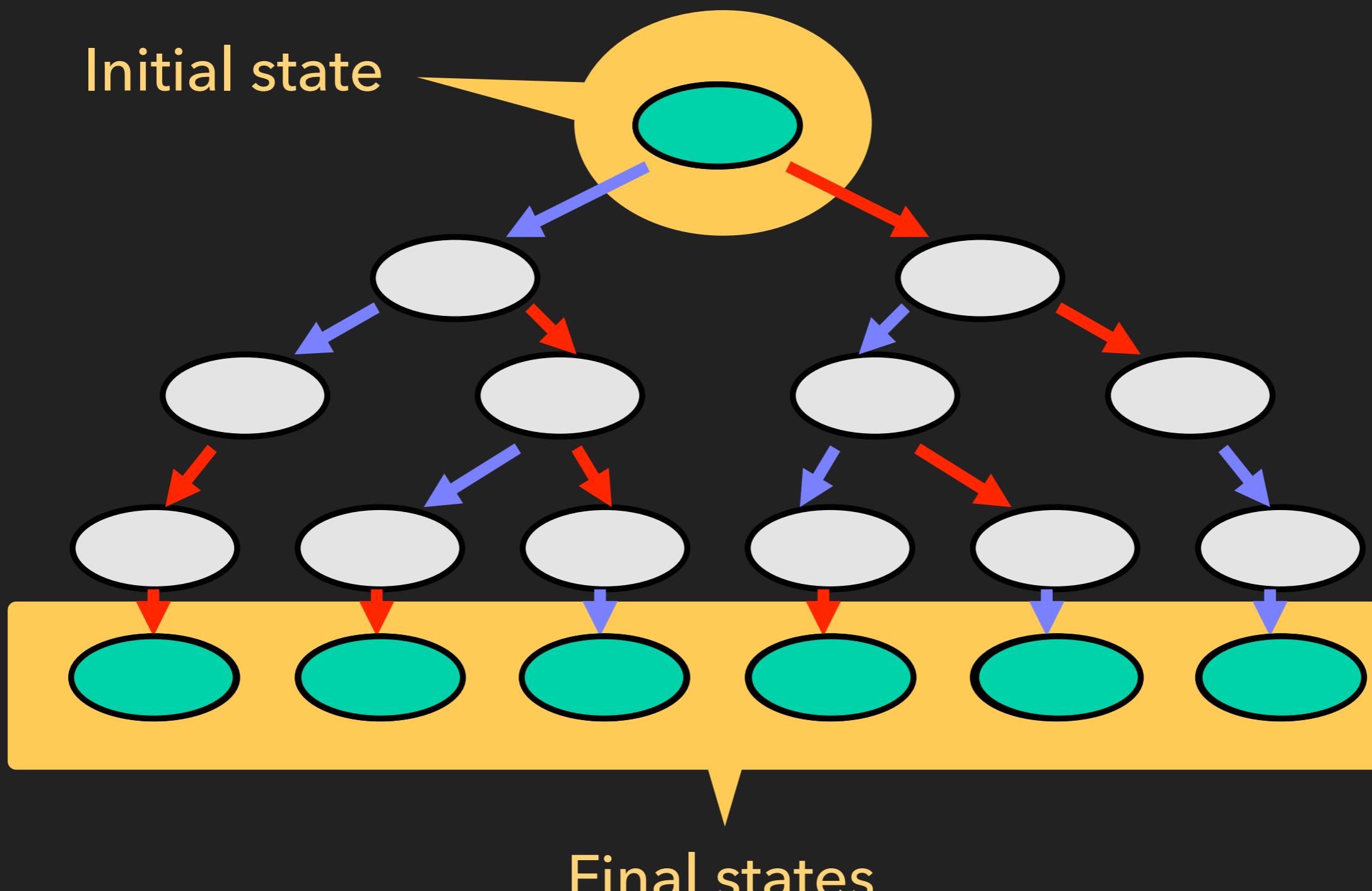
Quod  
Erat

Demonstrandum

Enough to consider binary consensus and  $n = 2$

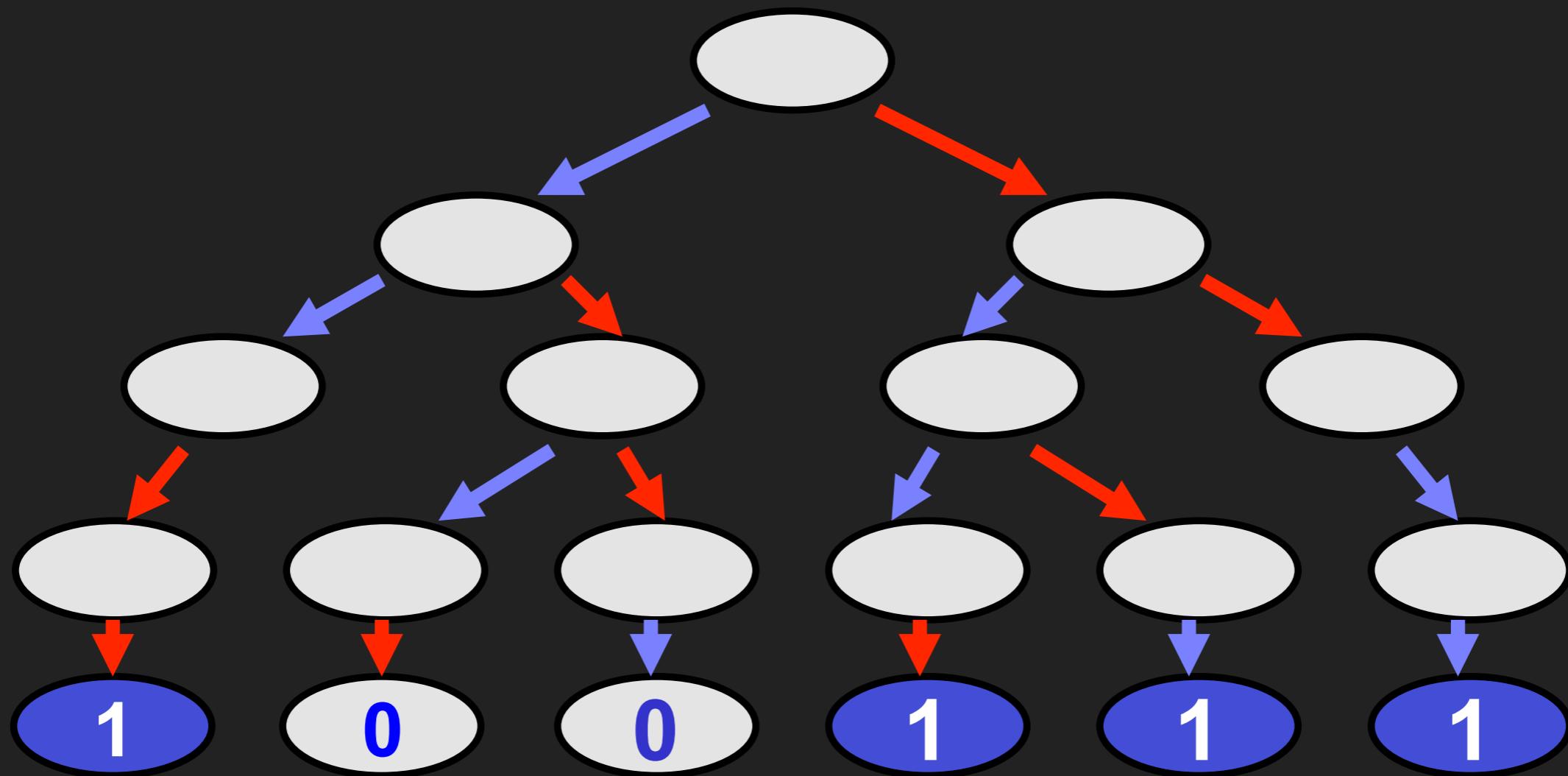


- ▶ Either A or B “moves”
- ▶ Moving means
  - ▶ Register read
  - ▶ Register write



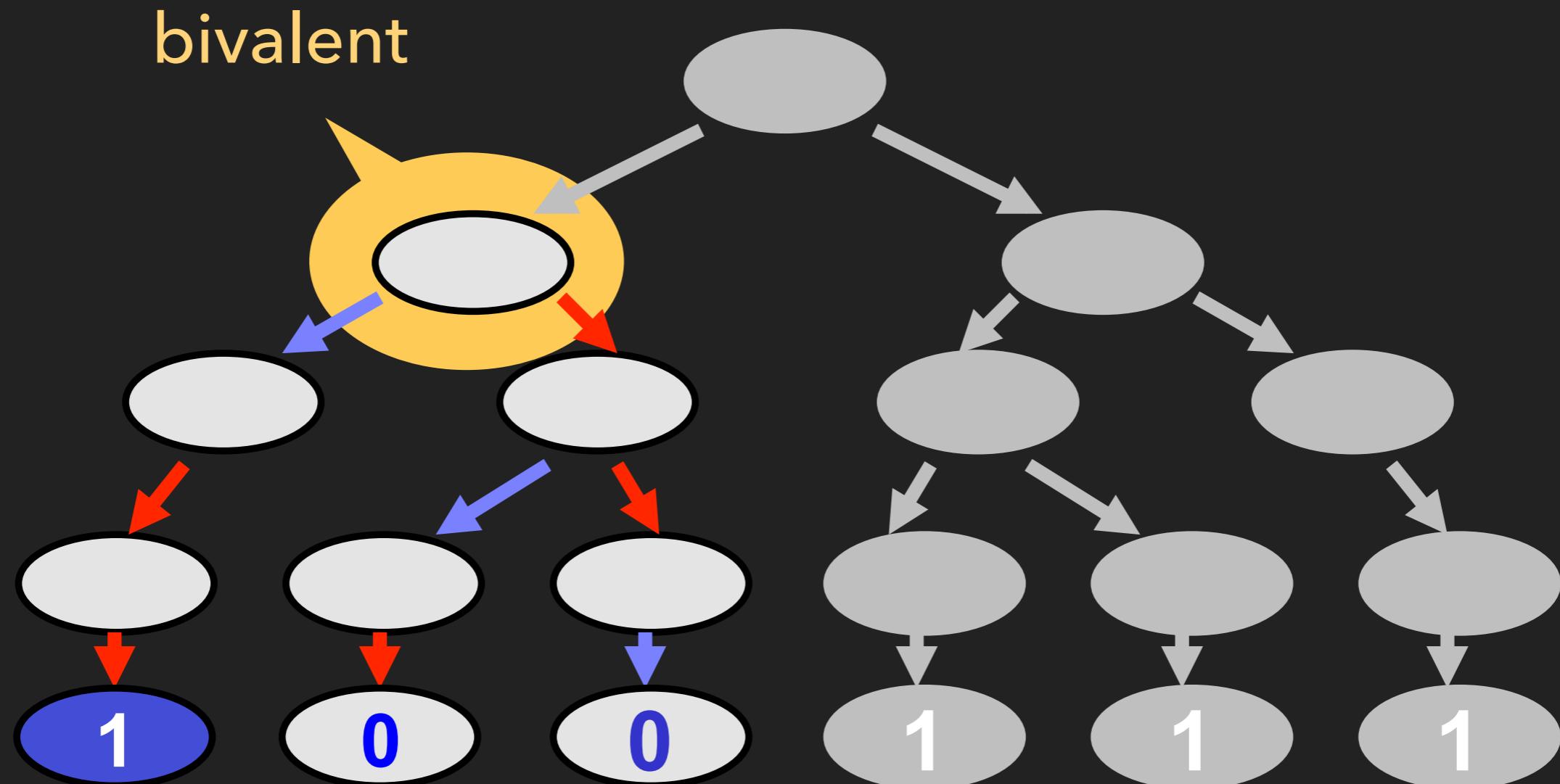
# DECISION VALUES

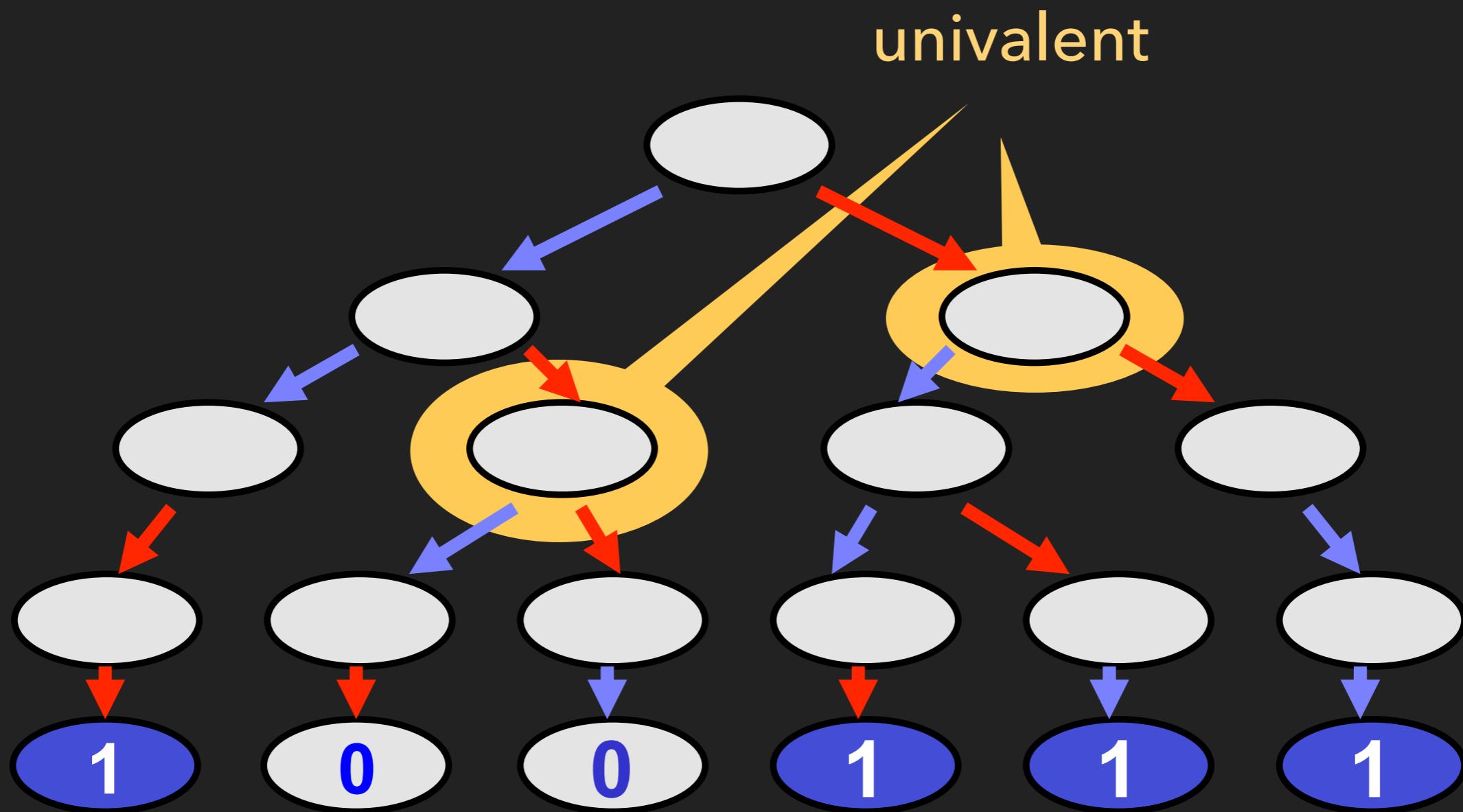
22

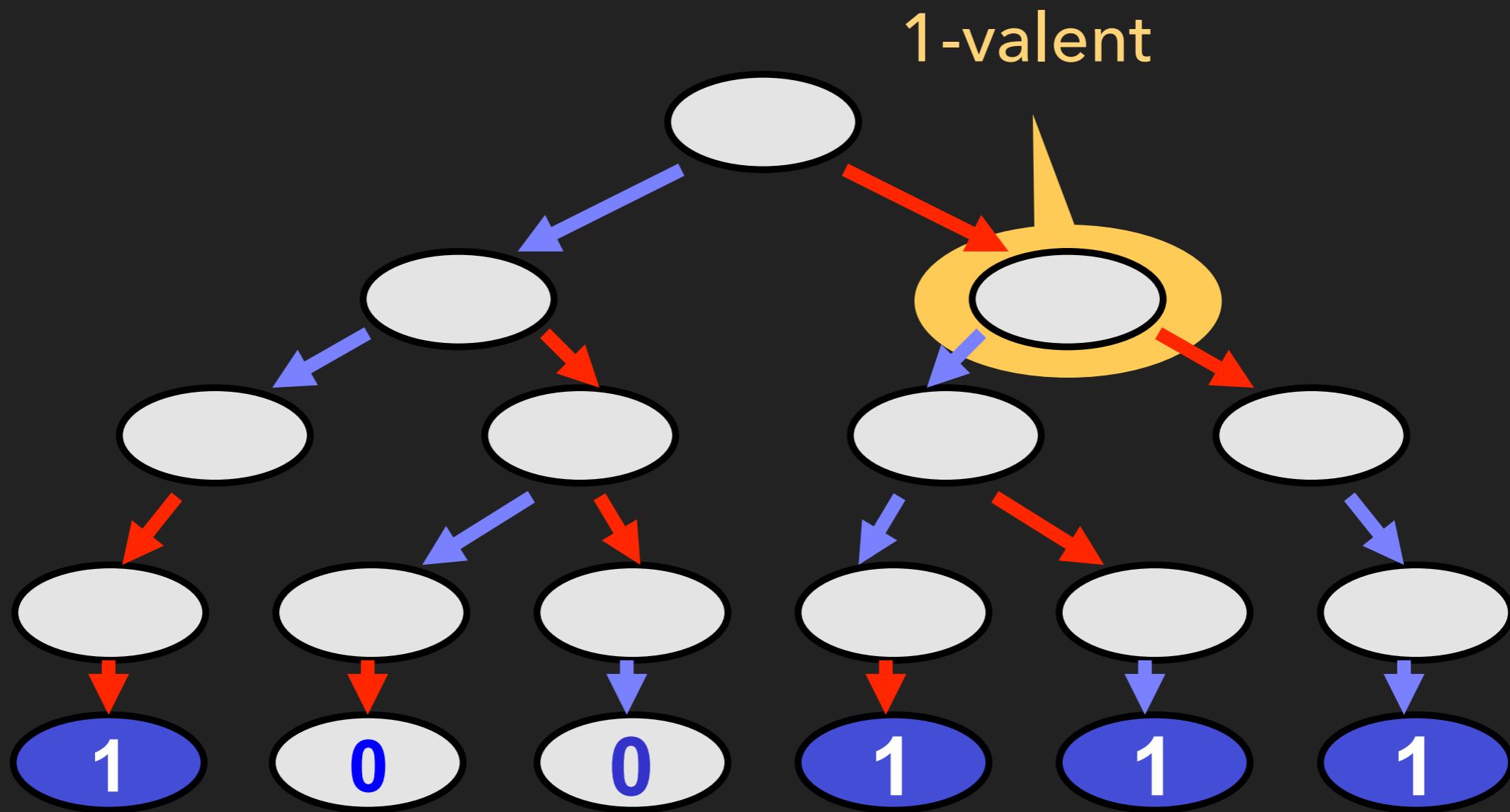


# BIVALENT : BOTH POSSIBLE

23



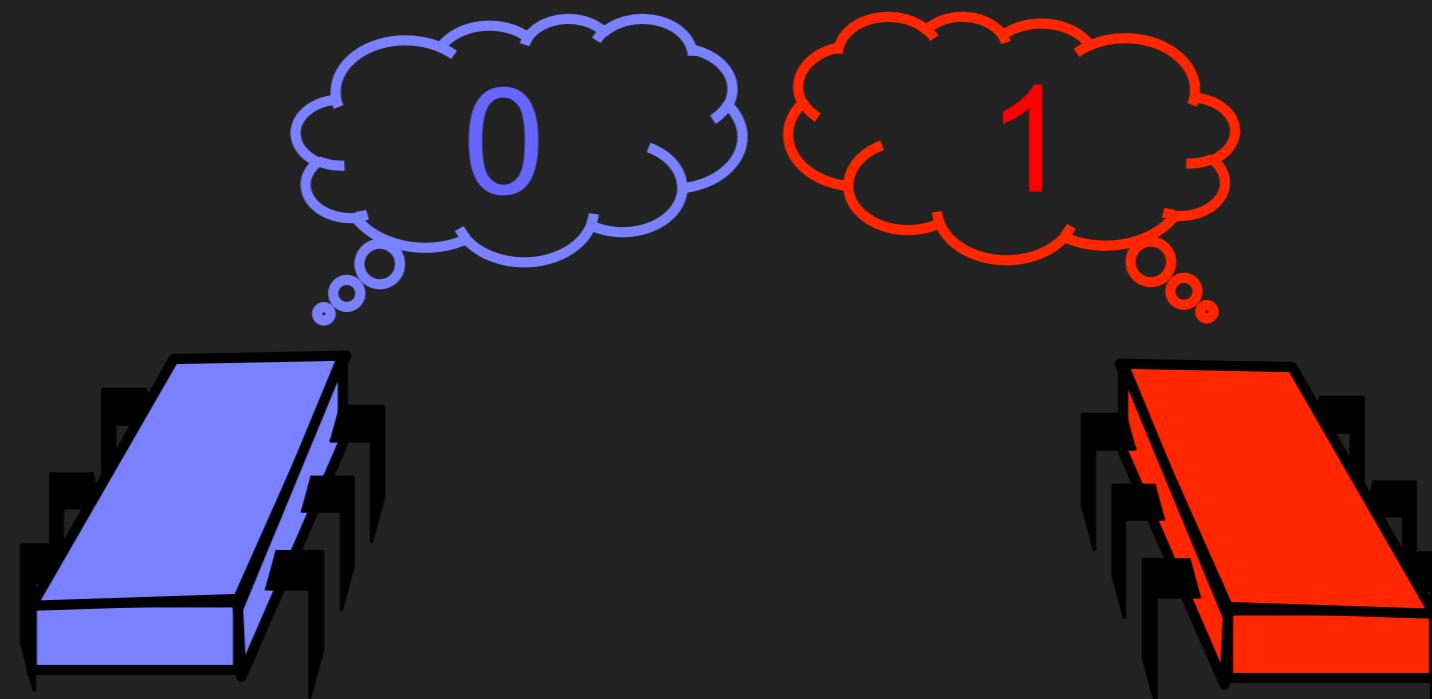




x-valent : x is the only decision value

- ▶ Any consensus protocol 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

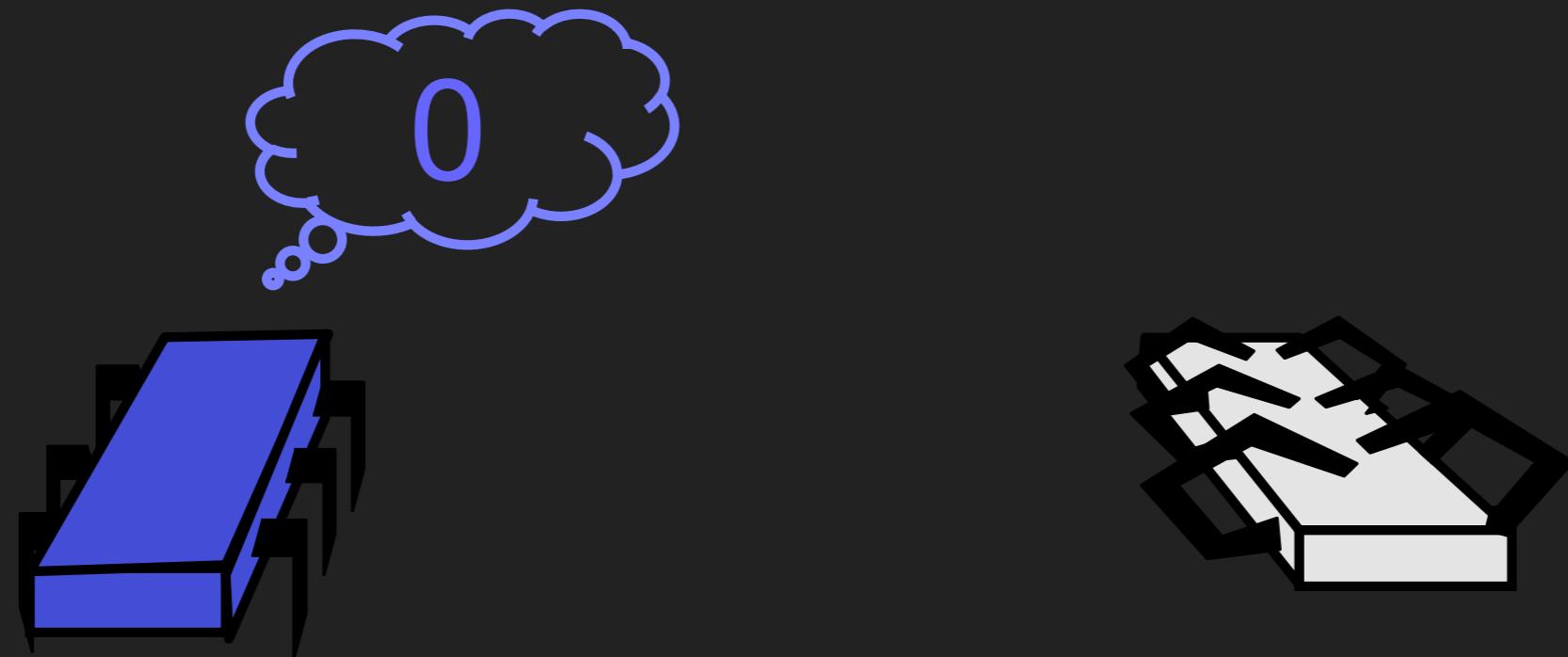
What if inputs differ?



# THERE EXISTS AN INITIAL BIVALENT STATE

---

Must decide on 0



In this solo execution by A

# THERE EXISTS AN INITIAL BIVALENT STATE

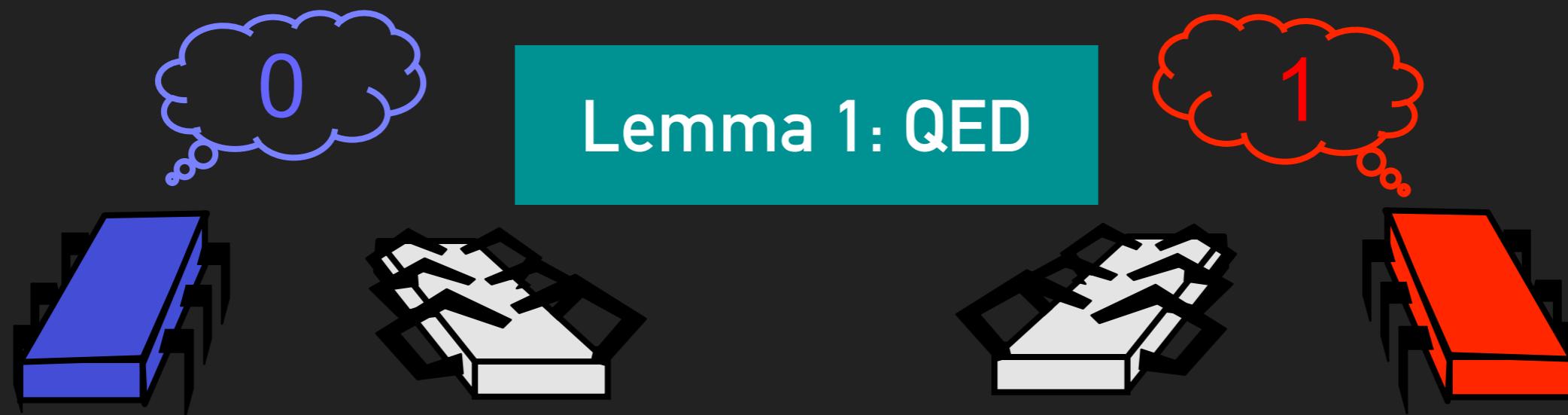
---

Must decide on 1



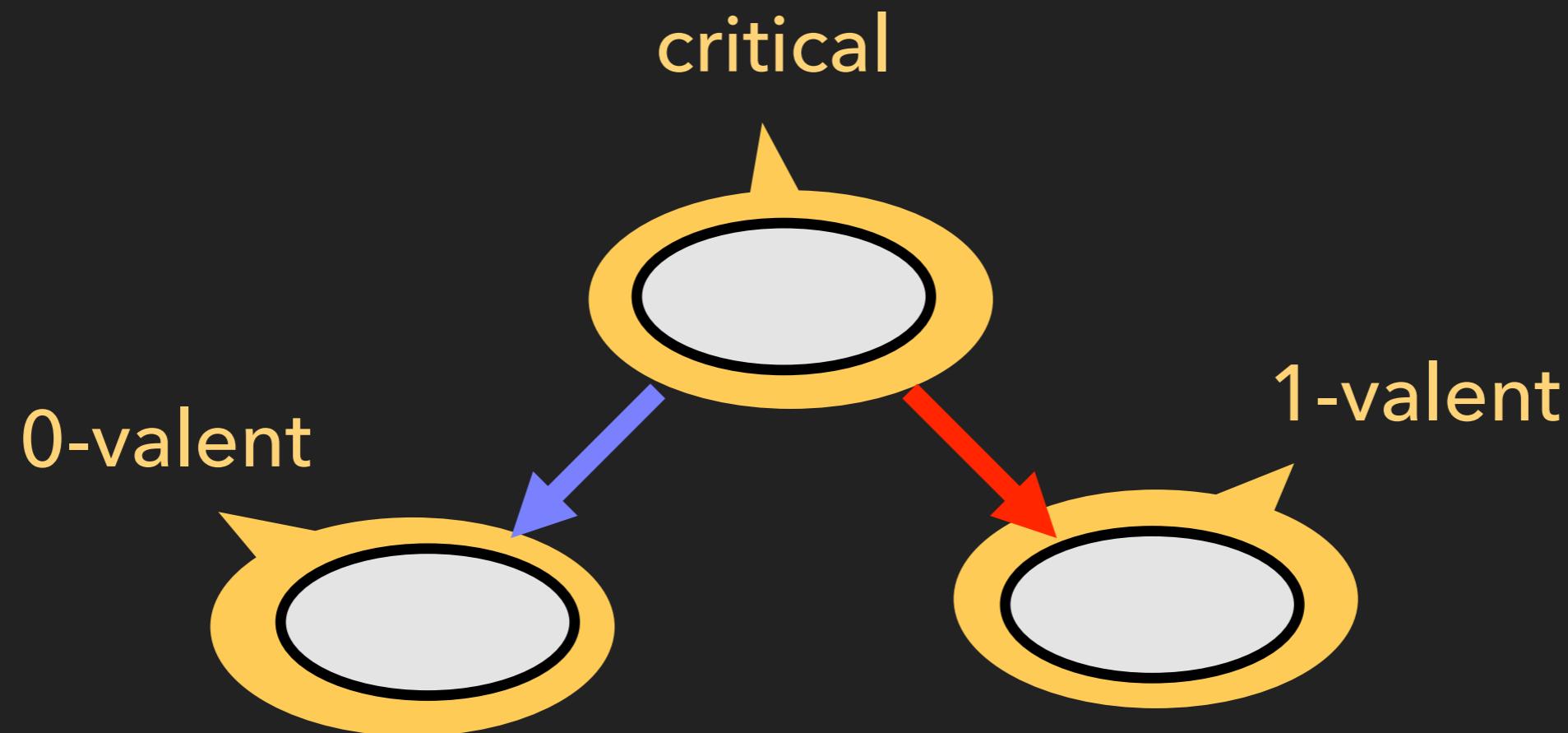
In this solo execution by B

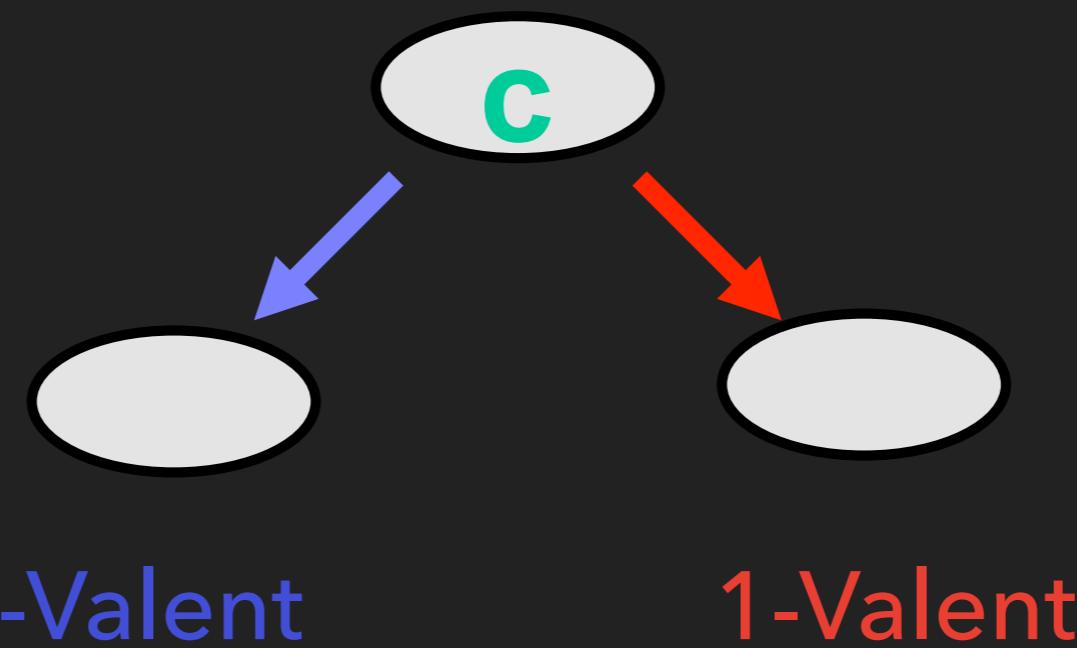
## Mixed initial state bivalent



Solo execution by A  
must decide 0

Solo execution by B  
must decide 1

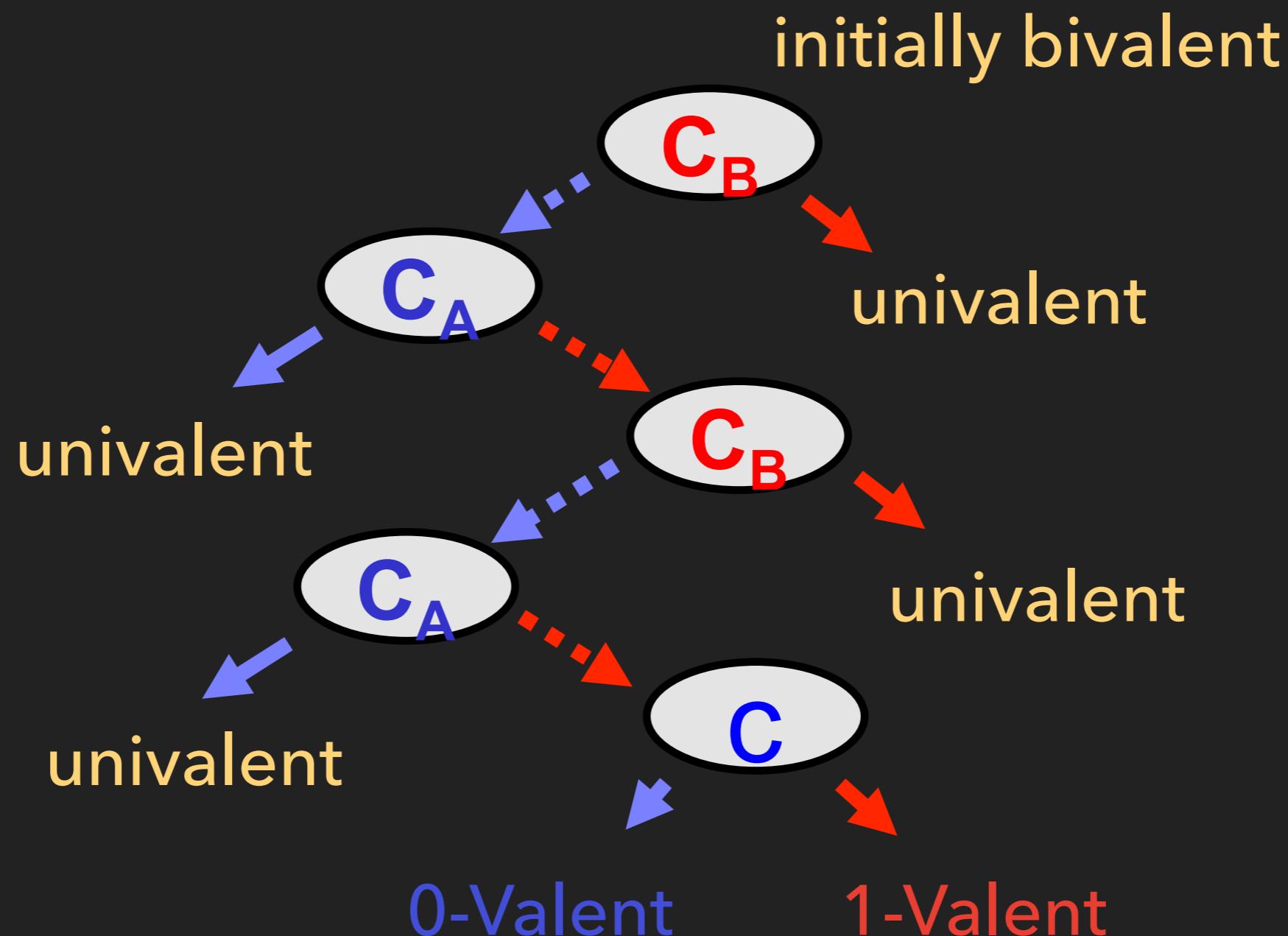




If A goes first protocol  
decides 0

If B goes first protocol  
decides 1

# WHY MUST THERE BE A CRITICAL STATE ?



# CRITICAL STATES

---

- ▶ Starting from a bivalent state
- ▶ The protocol can reach a critical state
  - ▶ otherwise we can stay bivalent forever
  - ▶ And the protocol does not terminate in finite steps

# CLOSING DEAL...

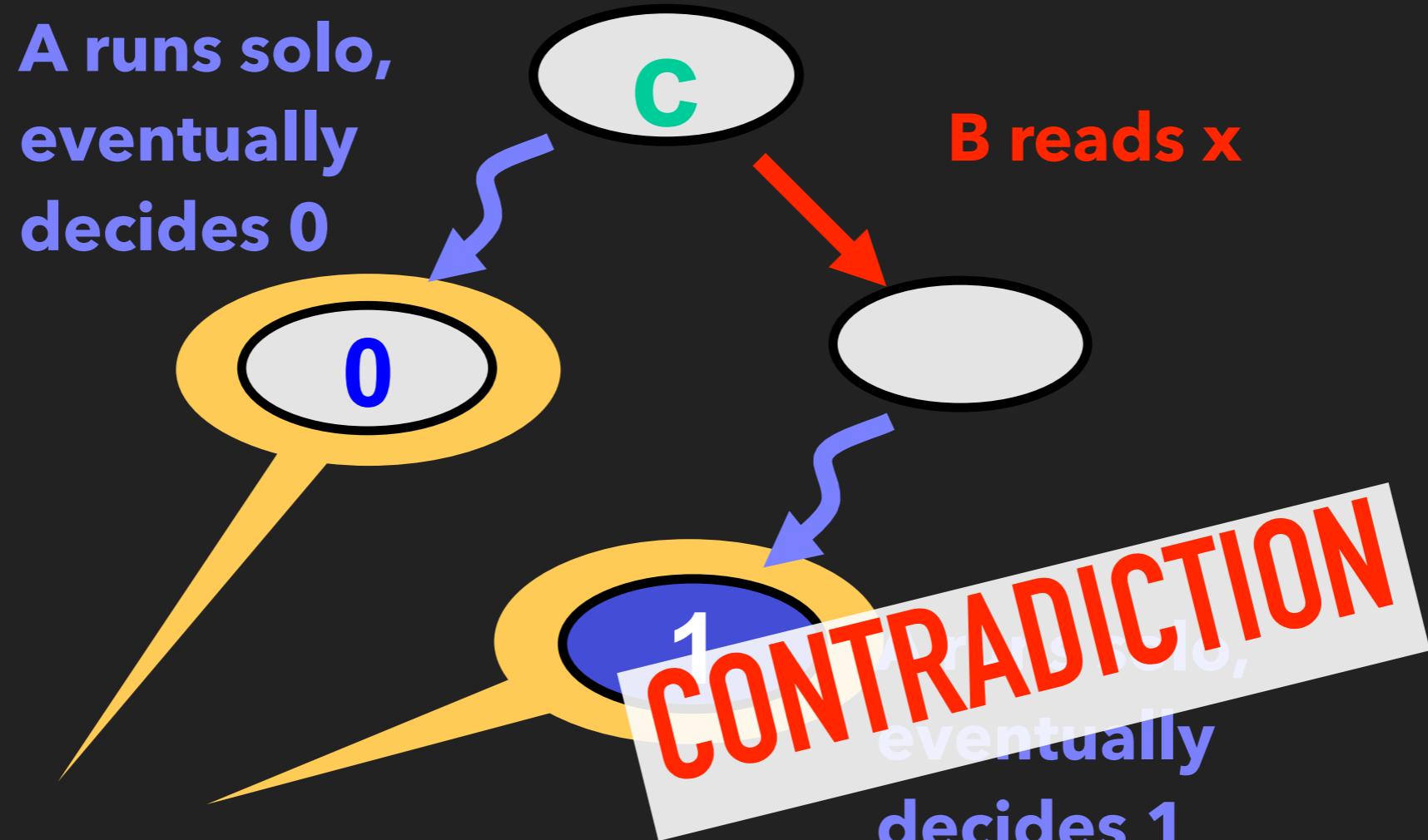
---

- ▶ Starting from a critical state
- ▶ Each processor fixes the outcome by
  - ▶ Reading or writing
  - ▶ Same or different registers
- ▶ Leading to a 0 or 1 decision
- ▶ And a contradiction

# POSSIBLE INTERACTIONS

36

		A reads x	A reads y	x.write	y.write
		x.read	y.read		
x.read	x.read	?	?	?	?
y.read	y.read	?	?	?	?
x.write	x.write	?	?	?	?
y.write	y.write	?	?	?	?



States look same to A

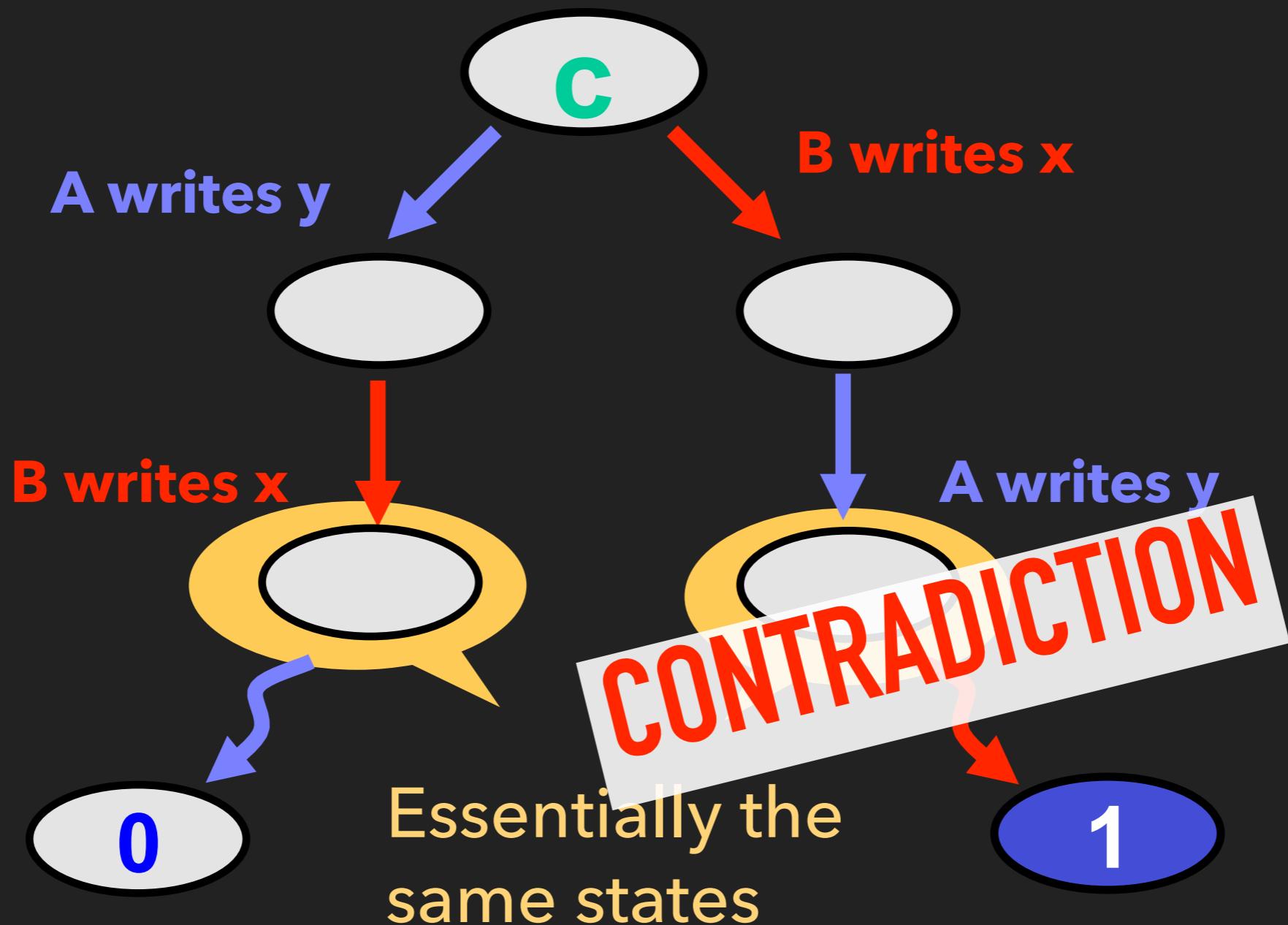
# POSSIBLE INTERACTIONS

38

	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

39



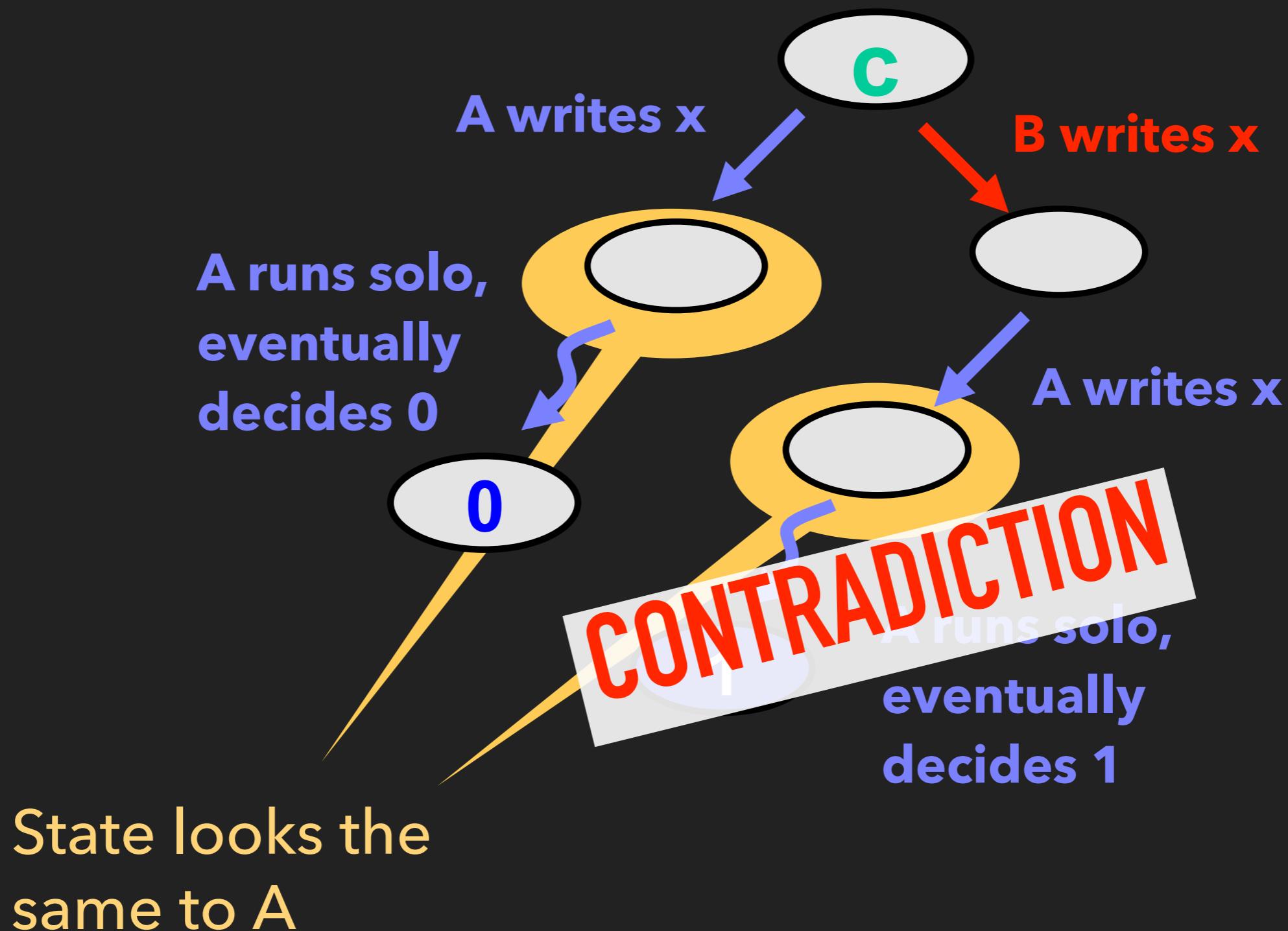
# POSSIBLE INTERACTIONS

40

	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

41



# POSSIBLE INTERACTIONS

42

	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

Q.E.D.

CONSENSUS\* IS IMPOSSIBLE IN A  
COMPLETELY ASYNCHRONOUS  
MESSAGE PASSING SYSTEMS WHERE  
EVEN A SINGLE PROCESSOR FAILS

[Fischer, Lynch, Paterson '85]

CONSENSUS\* IS IMPOSSIBLE IN A  
**COMPLETELY ASYNCHRONOUS**  
MESSAGE PASSING SYSTEMS WHERE  
EVEN A SINGLE PROCESSOR FAILS

[Fischer, Lynch, Paterson '85]

# WHAT ARE COMPLETELY ASYNCHRONOUS SYSTEMS

45

PROCESSOR ASYNCHRONY

MESSAGE ORDER ASYNCHRONY

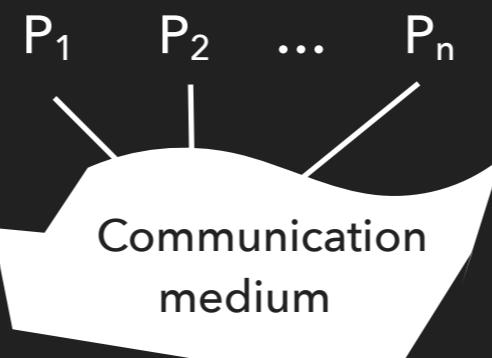
COMMUNICATION ASYNCHRONY

Are all three types of asynchrony needed simultaneously to obtain impossibility result?

NO

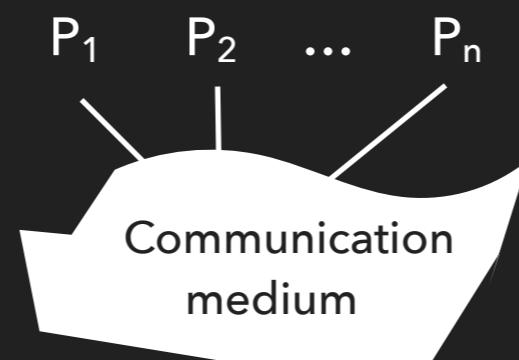
[Dolev, Dwork, Stockmeyer '87]

Minimal  
synchronism needed  
for distributed  
consensus



Processors can be  
**synchronous** or  
**asynchronous**

Communication  
delay can be  
**bounded** or  
**unbounded**

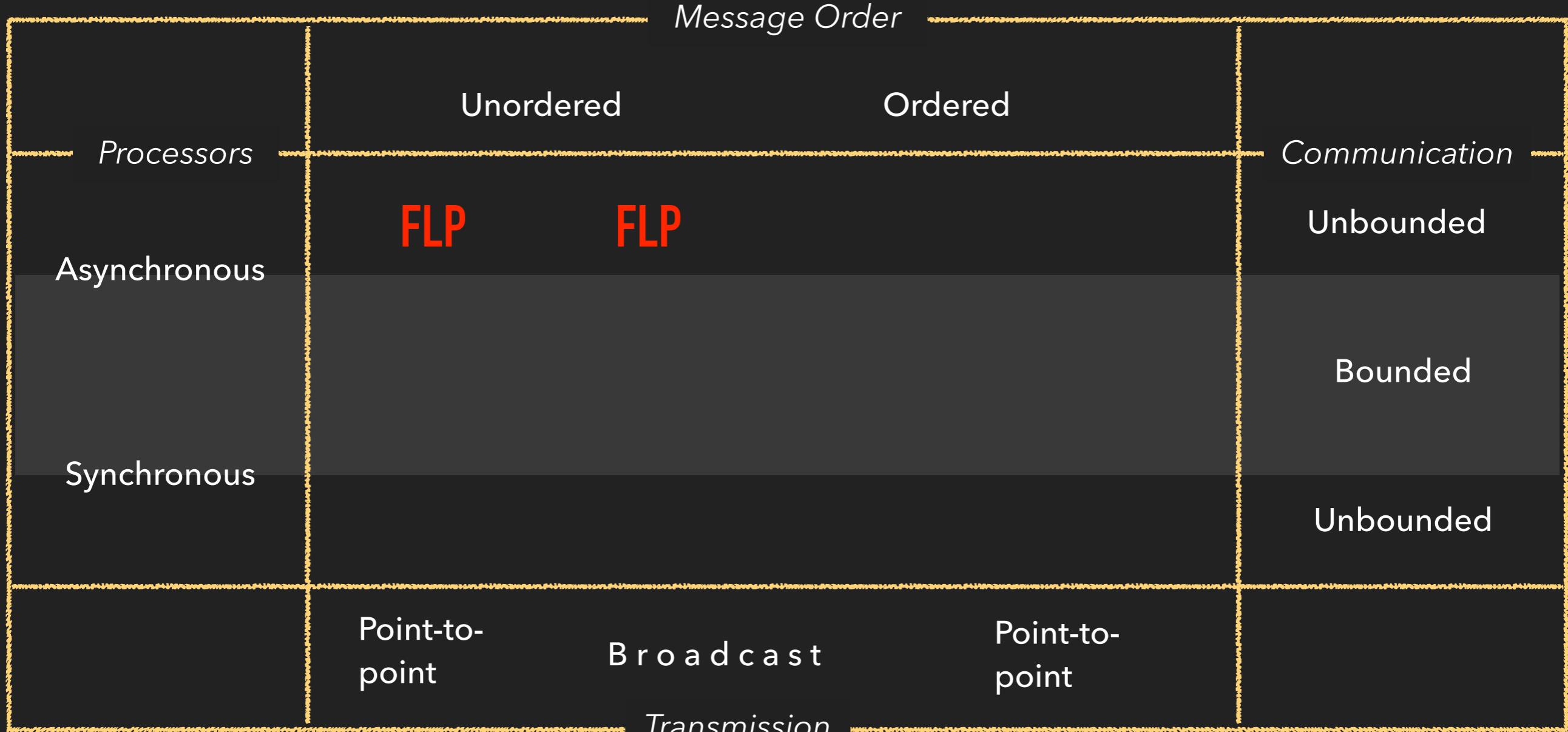


Messages can be  
**ordered** or  
**unordered**

Transmission  
mechanism can be  
**point-to-point** or  
**broadcast**

# WORLD OF (IM)POSSIBILITIES

48



16 possibilites

- ▶ 3 cases where N-resilient protocols exist

Case 1

Processors are synchronous and communication is bounded

Case 2

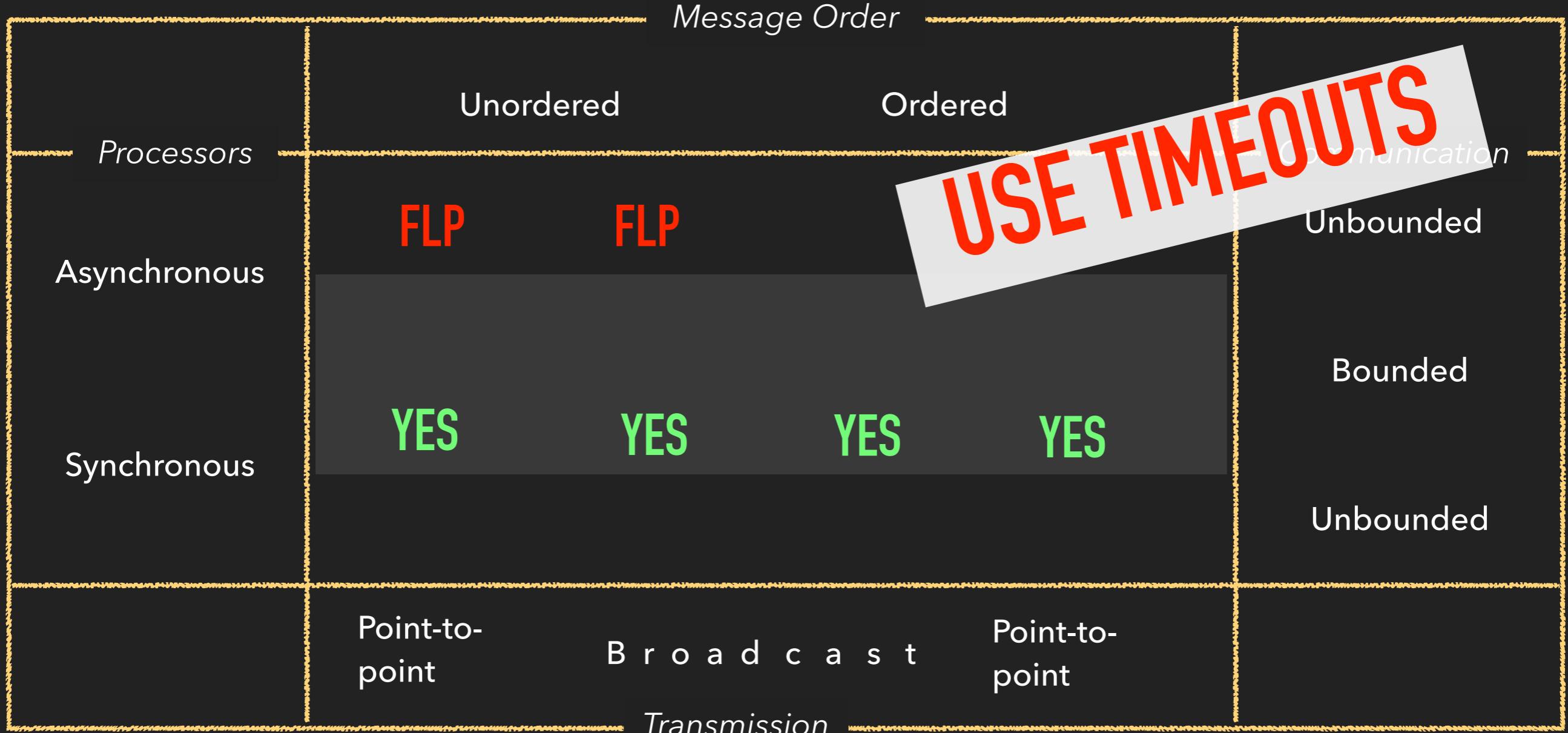
Messages are ordered and transmission mechanism is broadcast

Case 3

Processors are synchronous and messages are ordered

# WORLD OF (IM)POSSIBILITIES

50



Processors are synchronous and communication is bounded

DECIDE ON THE  
FIRST RECEIVED  
MESSAGE

		Message Order		Communication	
		Unordered	Ordered	Unbounded	Bounded
Processors	Asynchronous	FLP	FLP	YES	YES
	Synchronous	YES	YES	YES	YES
Transmission		Point-to-point	Broadcast	Point-to-point	Unbounded

Messages are ordered and transmission mechanism is broadcast

**EXPONENTIAL  
NUMBER OF  
MESSAGES**

		Message Order		Communication	
		Unordered	Ordered	Unbounded	Bounded
Processors	Asynchronous	FLP	FLP	YES	YES
	Synchronous	YES	YES	YES	YES
	Point-to-point			YES	YES
	Broadcast				
	Transmission			Point-to-point	

Processors are synchronous and messages are ordered

- ▶ 3 cases where N-resilient protocols exist
- ▶ No t-resilient protocol if system weakened for  $t = 1, 2$
- ▶ Favourable -> Unfavourable

Processors : synchronous / asynchronous

Communication : bounded / unbounded

Messages : ordered / unordered

Transmission : broadcast / point-to-point

# WORLD OF (IM)POSSIBILITIES

54

	Message Order		
Processors	Unordered	Ordered	Communication
Asynchronous	FLP	FLP	YES → NO
Synchronous	NO ↑ sync -> async	NO ↓ bdd -> u-bdd	YES YES YES
Point-to-point	YES	NO	YES YES
Broadcast			
Transmission			

# LAST LESSON BEFORE FIGHTING THE BEAST : PROBABILITY ROCKS<sup>55</sup>

---

- ▶ Probabilistic consensus possible in presence of faults
  - ▶ do not terminate wlp
    - ▶ if terminate, right result
  - ▶ decisions not consistent wlp

# PLOTTING A BYZANTINE AGREEMENT

---

# BYZANTINE GENERALS PROBLEM

[ Lamport, Shostak,  
Pease '82 ]

57



# BYZANTINE GENERALS PROBLEM

58



Attack

[A, R, A, R, A]



Retreat

[A, R, A, R, A]



Attack

[A, R, A, R, A]



Retreat

[A, R, A, R, A]



Attack

[A, R, A, R, A]

Choose the majority vote

Reach consensus



# BYZANTINE GENERALS PROBLEM

59



Attack

~~[A, R, A, R, A]~~

Do not care what bz decides



Retreat

[A, R, A, R, A]



Attack

[A, R, A, R, A]



Retreat

[A, R, A, R, A]



Attack

[A, R, A, R, A]



# BYZANTINE GENERALS PROBLEM

60



Attack

~~[A, R, A, R, A]~~

Do not care what bz decides



Retreat

[R, R, A, R, A]

Retreat

Consensus violated



Attack

[R, R, A, R, A]

Retreat



Retreat

[A, R, A, R, A]

Agree



Attack

[A, R, A, R, A]

Agree



- ▶ Is there a protocol such that

## CONSISTENCY

All **honest** parties agree on same value

## VALIDITY

Agreed-upon value is input to some **honest** party

## TERMINATION

Each **honest** party decides in finite number of steps

# WORLD OF (IM)POSSIBILITIES

62

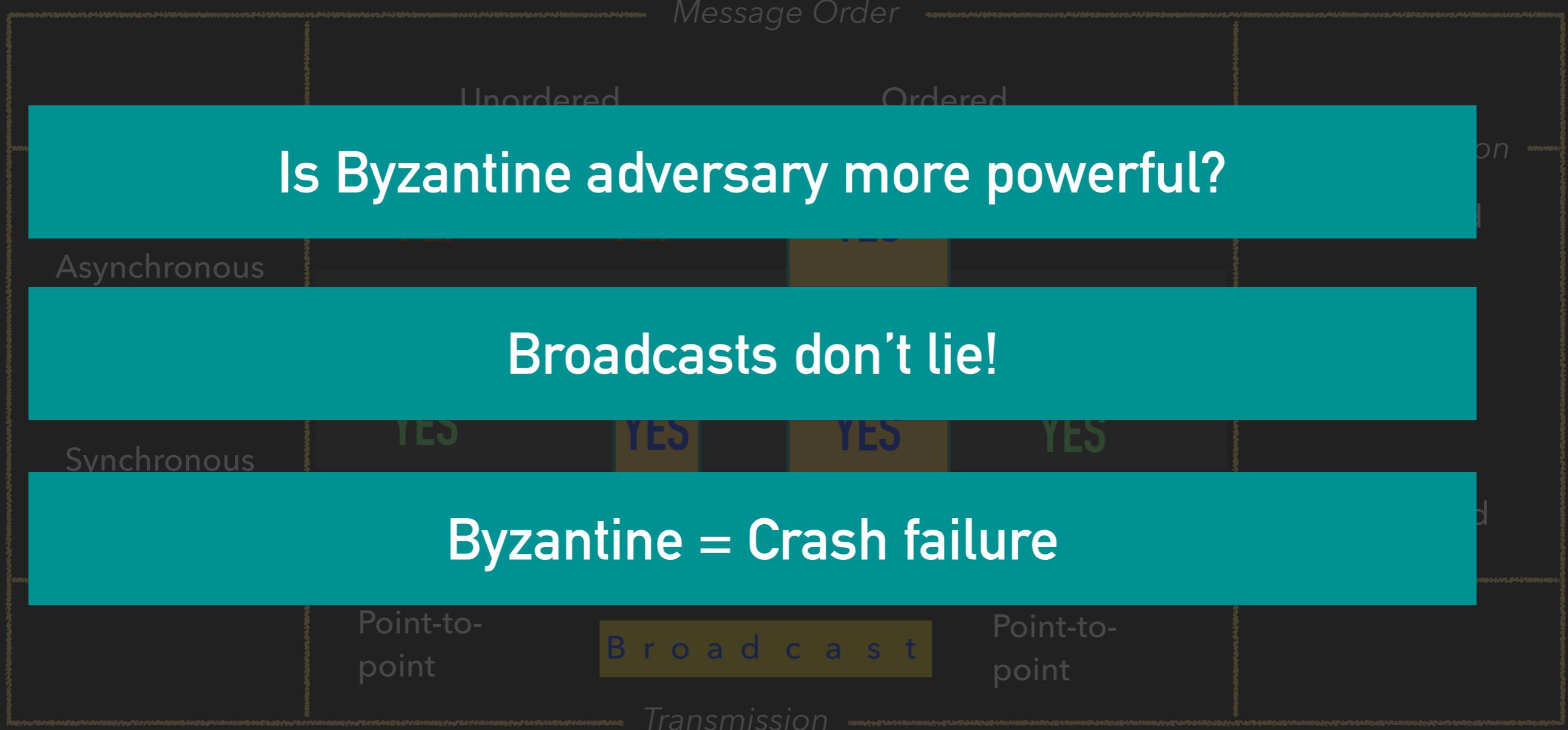
		Message Order					
		Unordered	Ordered				
Processors		FLP	FLP	YES	NO	Communication	
Asynchronous	FLP	NO	NO	YES	NO	Bounded	Unbounded
	YES	YES	YES	YES	YES		Unbounded
Synchronous		NO	NO	YES	YES		
Point-to-point		Broadcast		Point-to-point		Transmission	

Impossible in CF model => Impossible in Byzantine model

# WORLD OF (IM)POSSIBILITIES

63

		Message Order			
		Unordered	Ordered		
Processors		FLP	FLP	YES	NO
Asynchronous		NO	NO	YES	NO
Synchronous		YES	YES	YES	YES
	Point-to-point	NO	NO	YES	YES
		Broadcast		Point-to-point	
		Transmission		Communication	
		Unbounded		Bounded	
		Unbounded			



# WORLD OF (IM)POSSIBILITIES

		Message Order					
		Unordered		Ordered			
Processors		FLP	FLP	YES	NO	Communication	
Asynchronous		NO	NO	YES	NO	Unbounded	
Synchronous		YES	YES	YES	YES	Bounded	
		NO	NO	YES	YES	Unbounded	
Point-to-point		Broadcast				Point-to-point	
		Transmission					

BYZANTINE CONSENSUS IS  
IMPOSSIBLE WHEN NO  
AUTHENTICATION IS POSSIBLE IF  $T \geq 1/3 N$

[Lamport, Shostak, Pease '82]

Assume otherwise ...

Reason about the properties of any such protocol

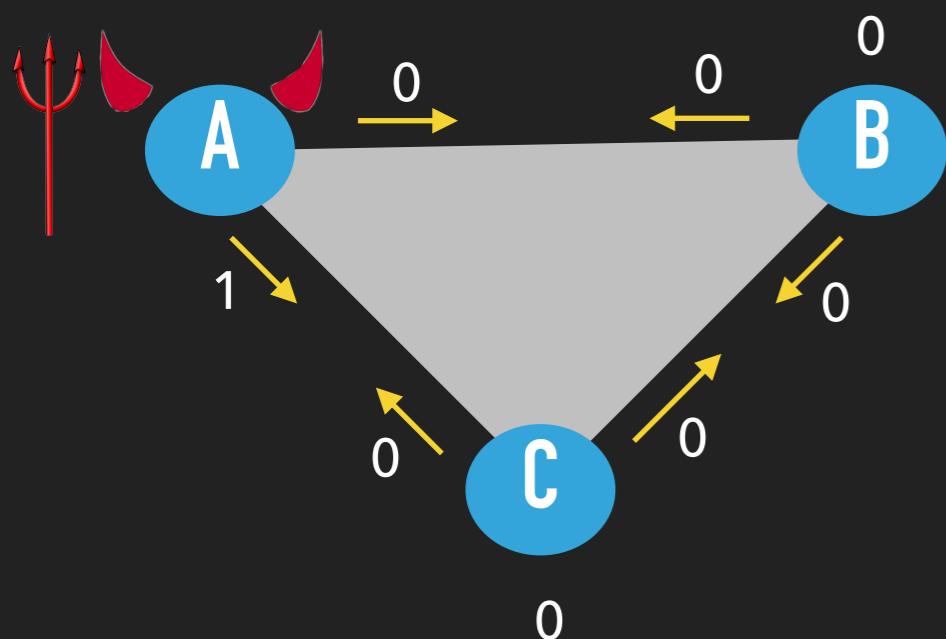
Derive a contradiction

Quod  
Erat  
Demonstrandum

Show for  $n = 3$  and  $t = 1$

# IMPOSSIBILITY PROOF (HIGH LEVEL)

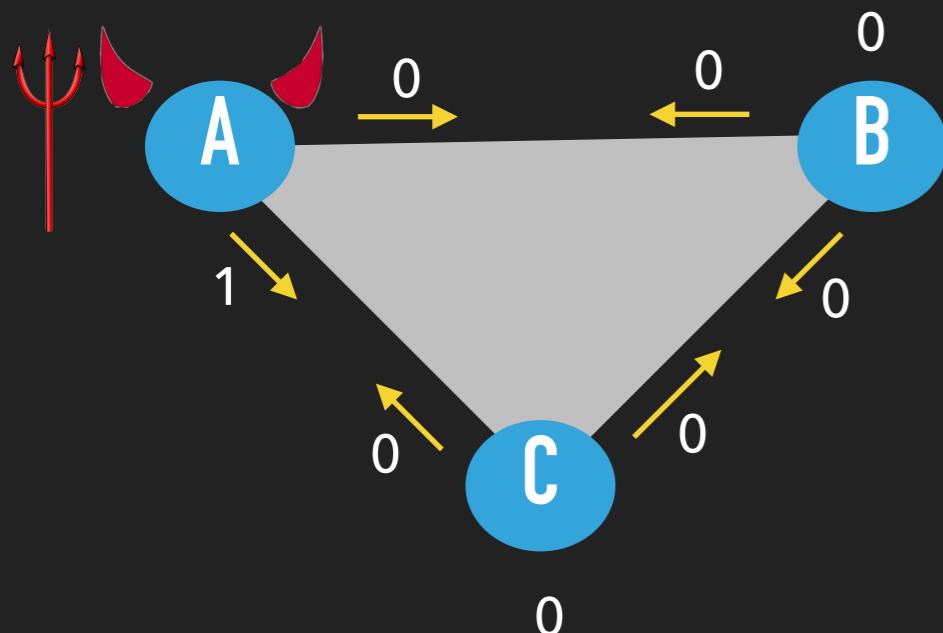
68



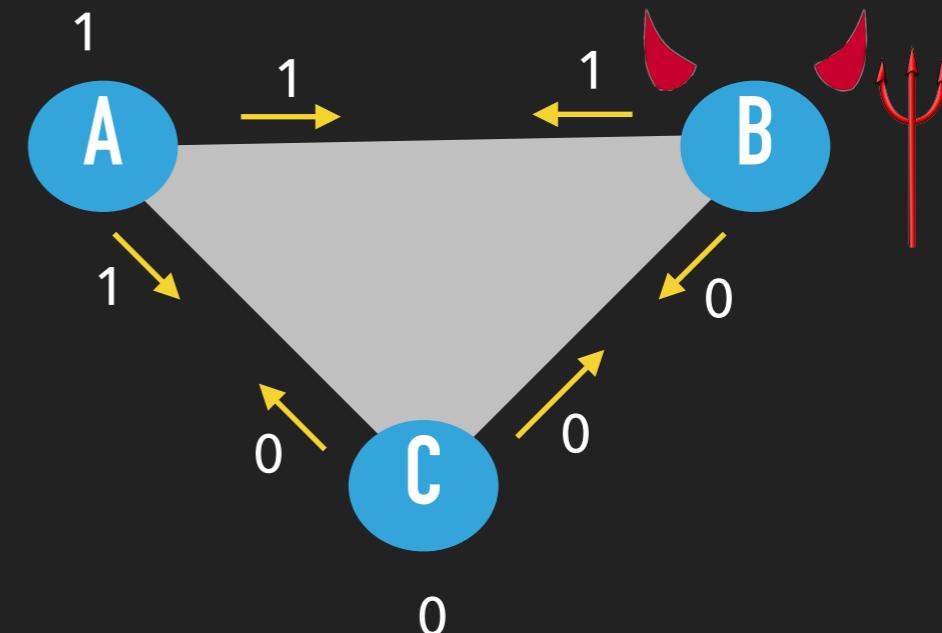
validity: Decide 0

# IMPOSSIBILITY PROOF (HIGH LEVEL)

69



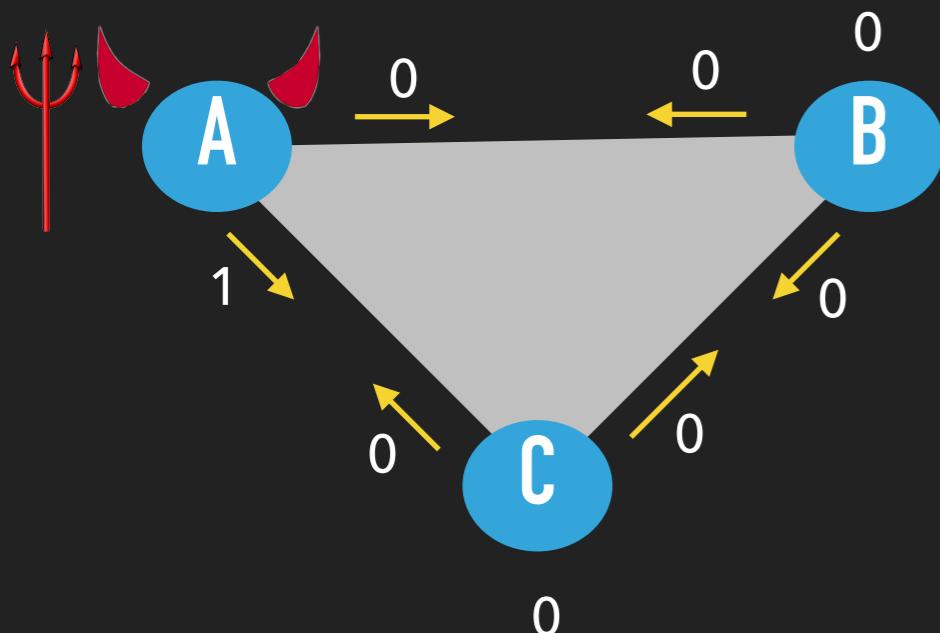
validity: Decide 0



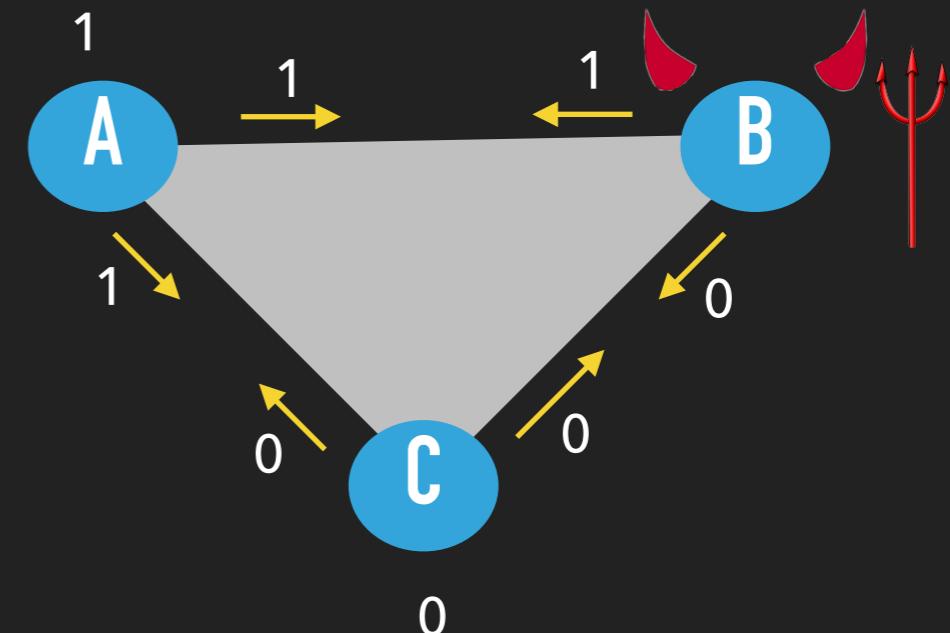
Decide 0 : View of C is  
indistinguishable from Case 1

# IMPOSSIBILITY PROOF (HIGH LEVEL)

70



validity: Decide 0

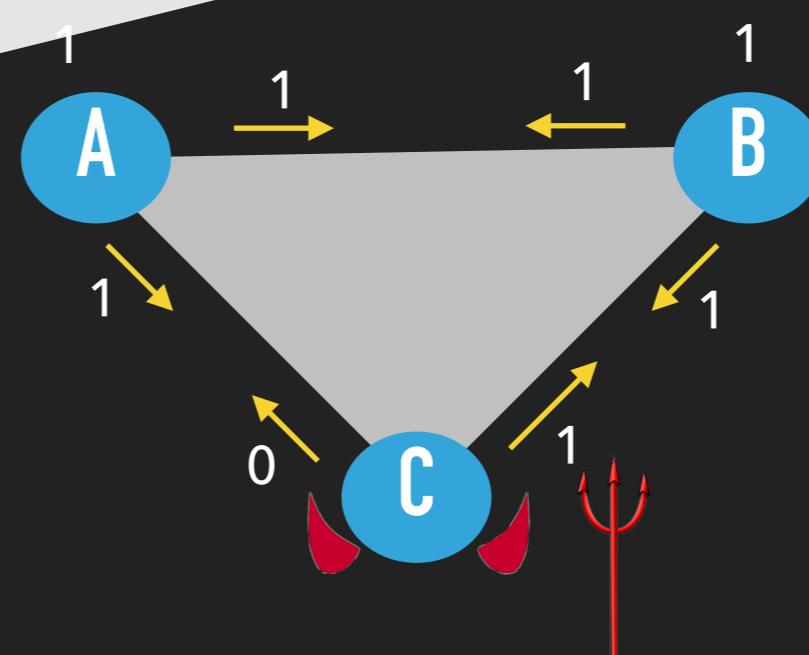


Decide 0 : View of C is  
indistinguishable from Case 1

**CONTRADICTION**

validity: Decide 1

A decides 0: View  
indistinguishable from Case 2

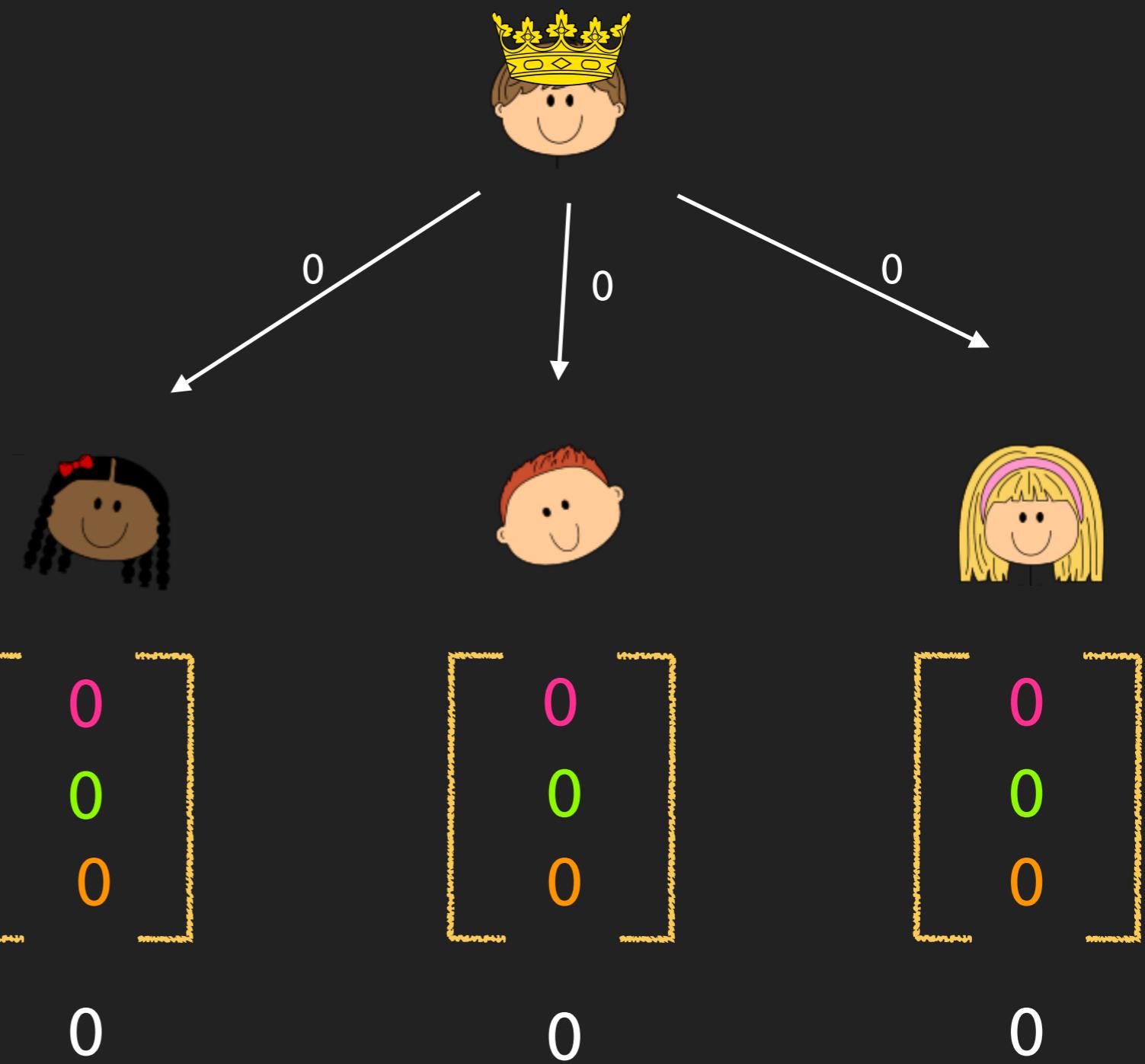


# WORLD OF (IM)POSSIBILITIES

71

		Message Order					
		Unordered		Ordered			
Processors		FLP	FLP	YES	NO	Communication	
Asynchronous		NO	NO	YES	NO	Unbounded	
Synchronous		YES*	YES	YES	YES*	Bounded	
		NO	NO	YES	YES*	Unbounded	
Point-to-point		Broadcast				Point-to-point	
		Transmission					

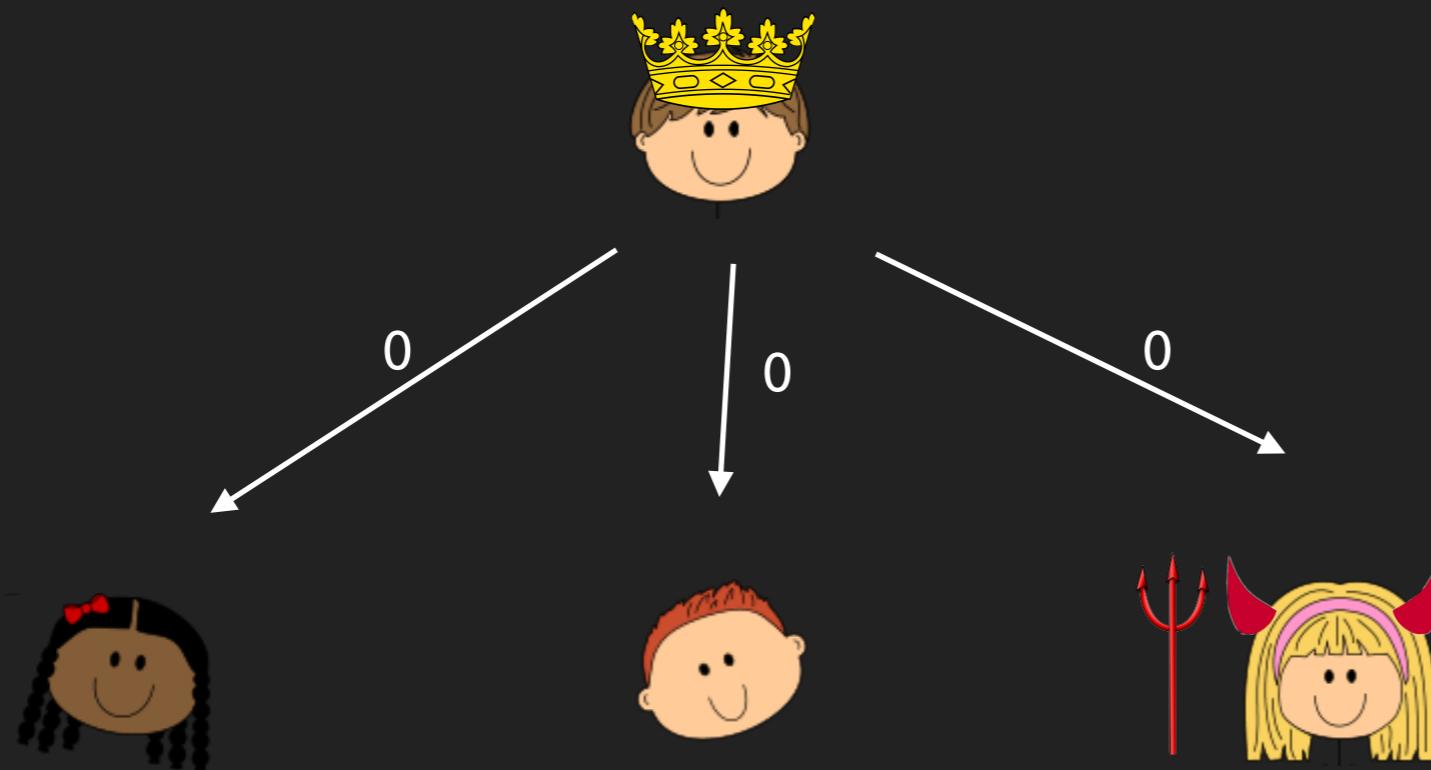
- ▶ Case n = 4 , t = 1
- ▶ Choose a leader
- ▶ Listen to what he says
- ▶ Talk to others
  - ▶ See what they say about the leader said
  - ▶ Decide on the majority value



# ALGORITHM IDEA WHEN $T < 1/3 N$

73

Case 1



CONSISTENCY

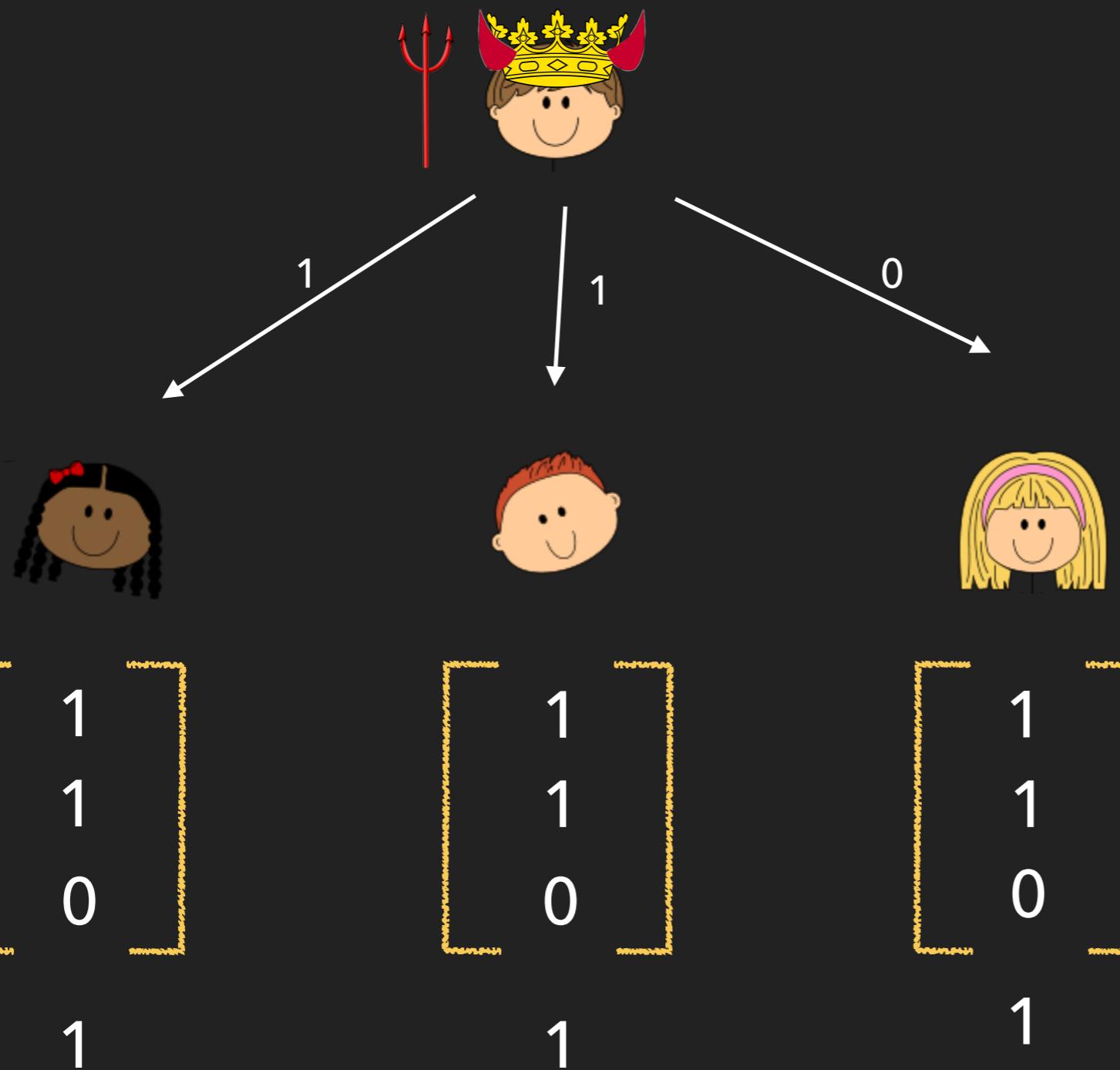
VALIDITY

[0 0 ?]	[0 0 ?]	[? ? 0]
0	0	

# ALGORITHM IDEA WHEN T < 1/3 N

74

Case 2



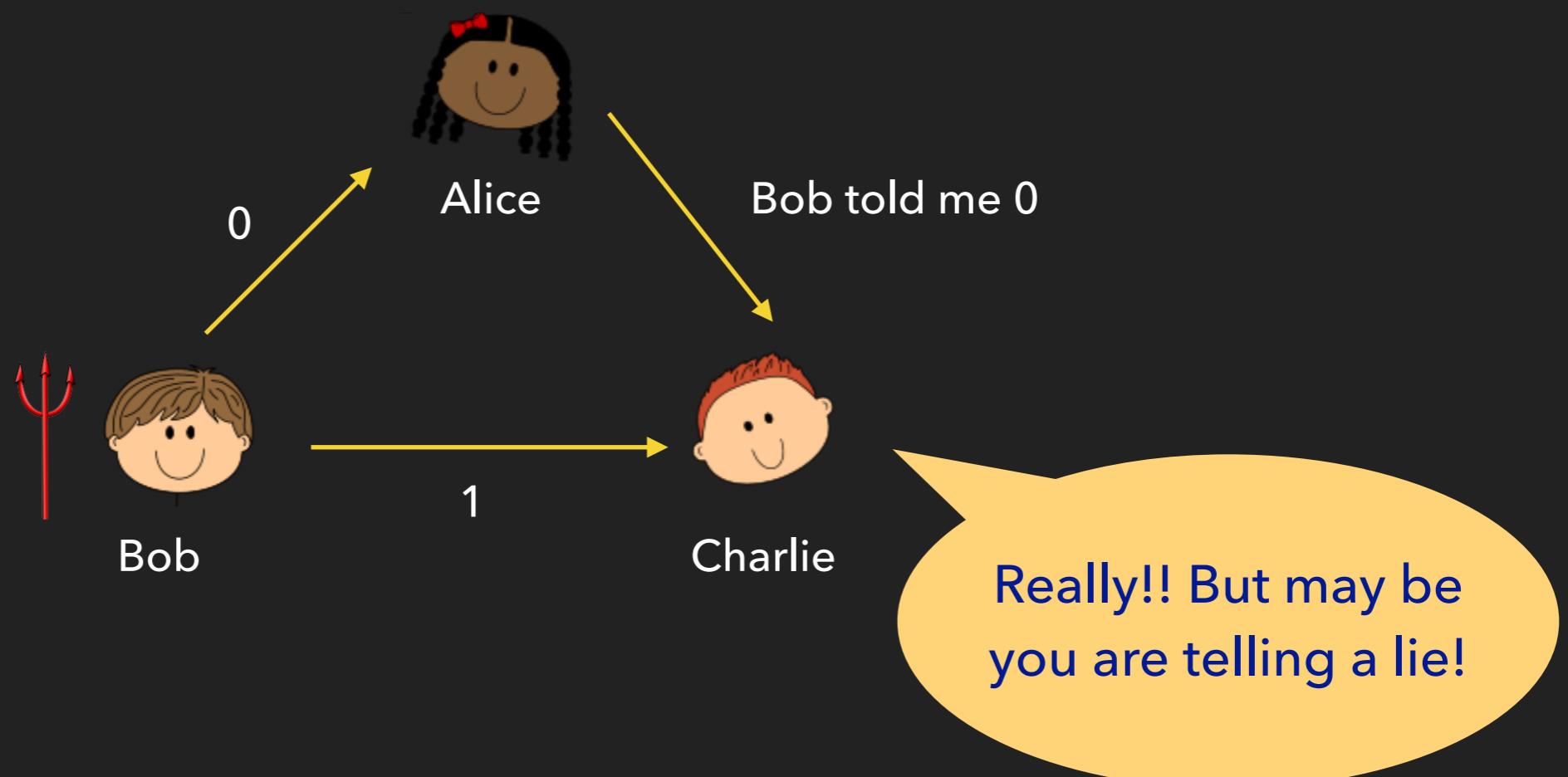
CONSISTENCY

VALIDITY

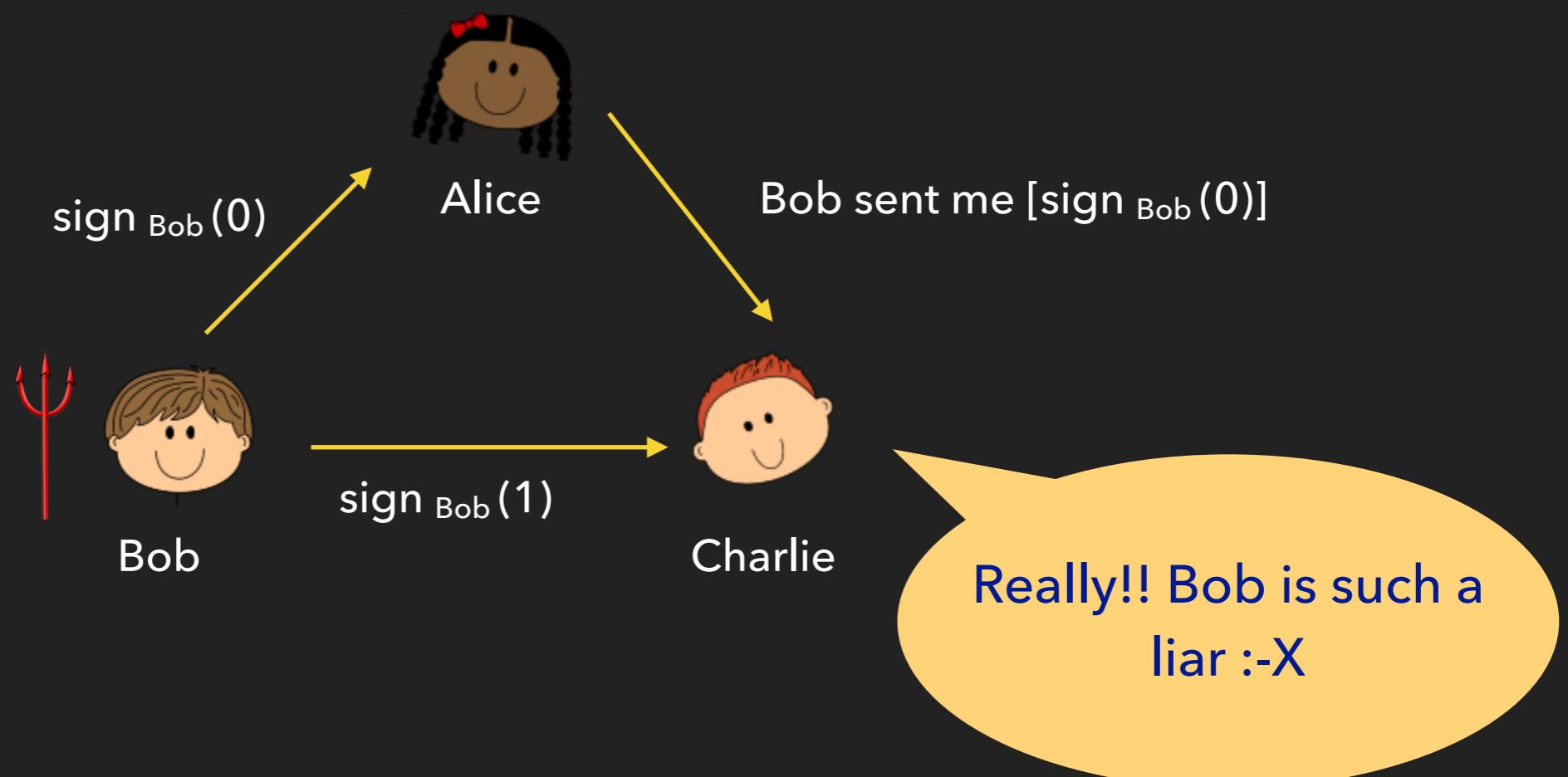
BYZANTINE CONSENSUS IS  
IMPOSSIBLE WHEN NO  
**AUTHENTICATION** IS POSSIBLE IF T  
 $>= 1/3 N$

[Lamport, Shostak, Pease '82]

- ▶ Malicious can send conflicting messages to other two



- ▶ Malicious can send conflicting messages to other two



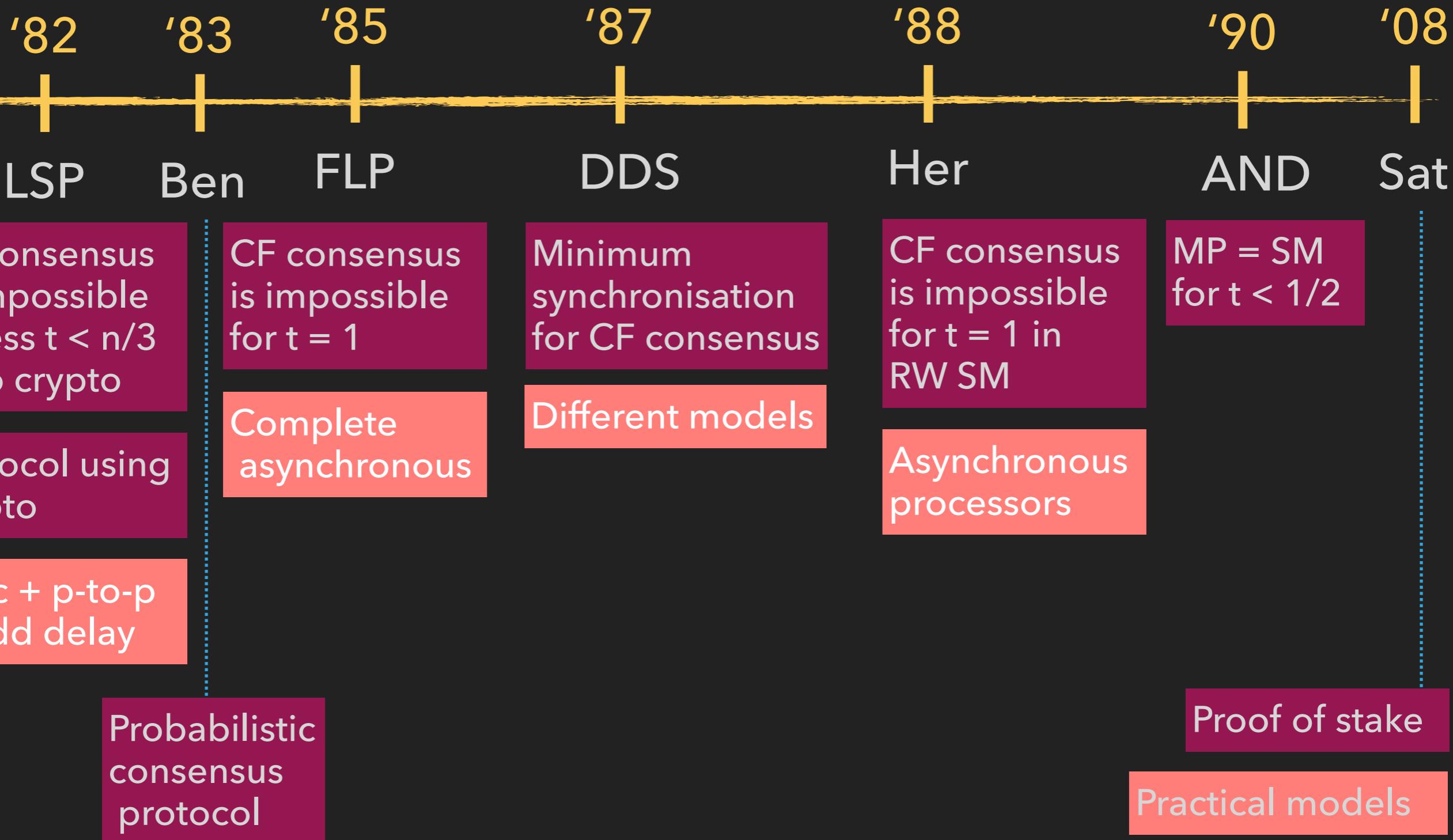
# SIGN ON THE DOTTED LINE [Lamport, Shostak, Pease '82]

---

- ▶ Use authenticated messages
- ▶ Give byzantine consensus protocol
  - ▶ Remember the model (sync processors + bdd delay)
- ▶ Use cryptographic techniques to sign

# THE ADVANCES IN CONSENSUS : SANDS OF TIME

79



Crash failure consensus is hard

Byzantine is even harder

Crypto is cool

that's why we have been doing it

Distributed + crypto is even cooler

that's why we will be doing it

Questions?

# REFERENCES

---

- ▶ Crash failure impossibility proof from lecture slides of Maurice Herlihy and Nir Shavit
- ▶ Many Faces of Consensus in Distributed Systems by John Turek and Dennis Shasha

# OVERVIEW

---