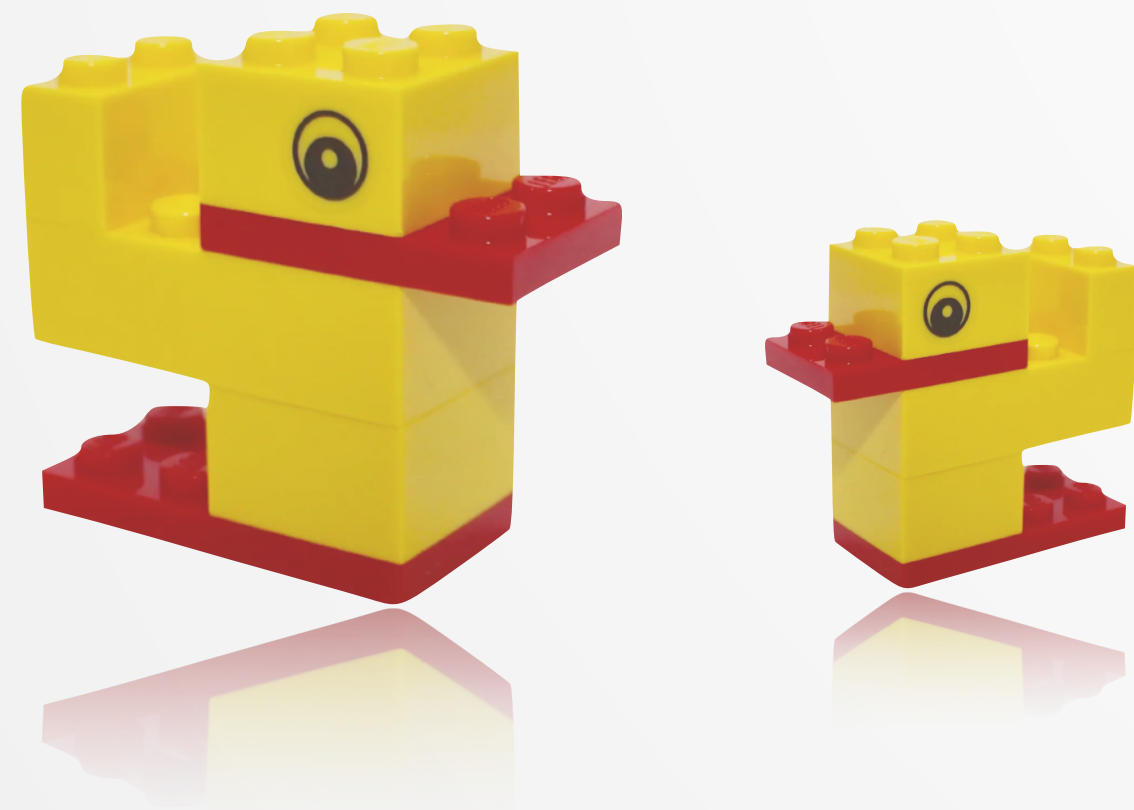**Advanced Course**

# Distributed Systems
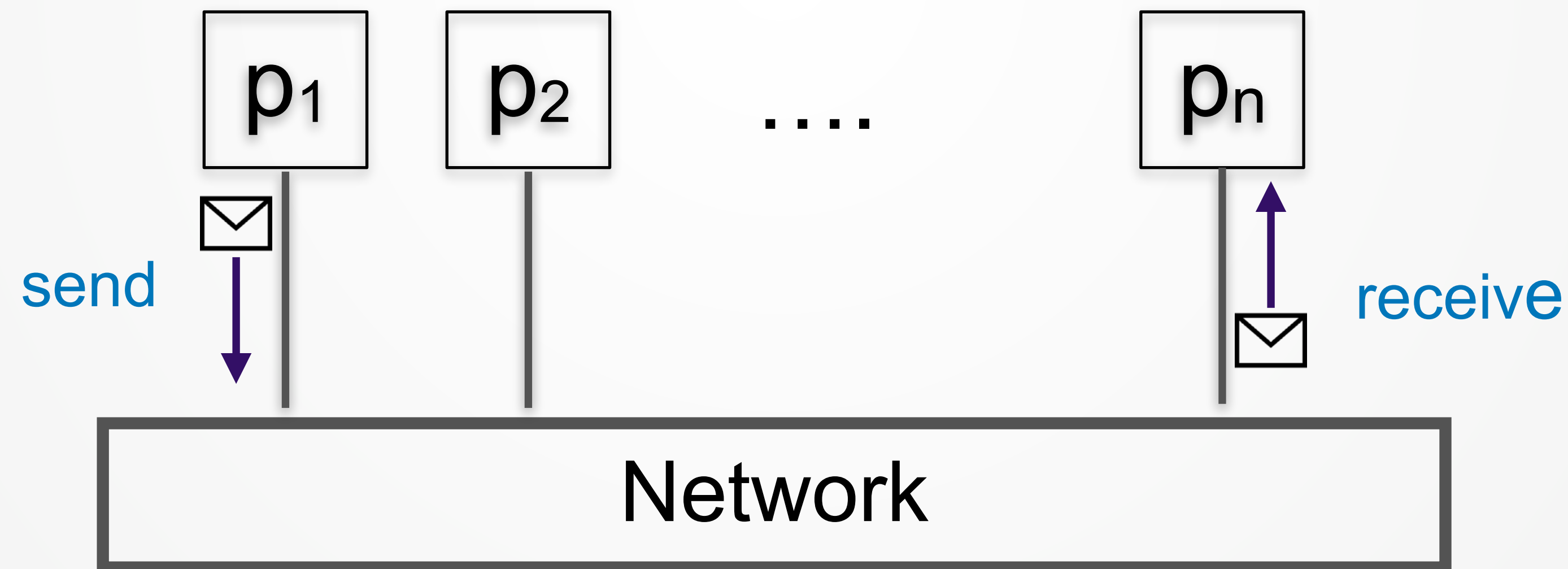
## Introduction to Distributed Systems

Paris Carbone

# COURSE TOPICS

▶ Intro to Distributed Systems

▶ Basic Abstractions and Failure Detectors

▶ Reliable and Causal Order Broadcast

▶ Distributed Shared Memory

▶ Consensus (Paxos, Raft, etc.)

▶ Dynamic Reconfiguration

▶ Time Abstractions and Interval Clocks (Spanner etc.)

▶ Consistent Snapshotting (Stream Data Management)

▶ Distributed ACID Transactions (Cloud DBs)
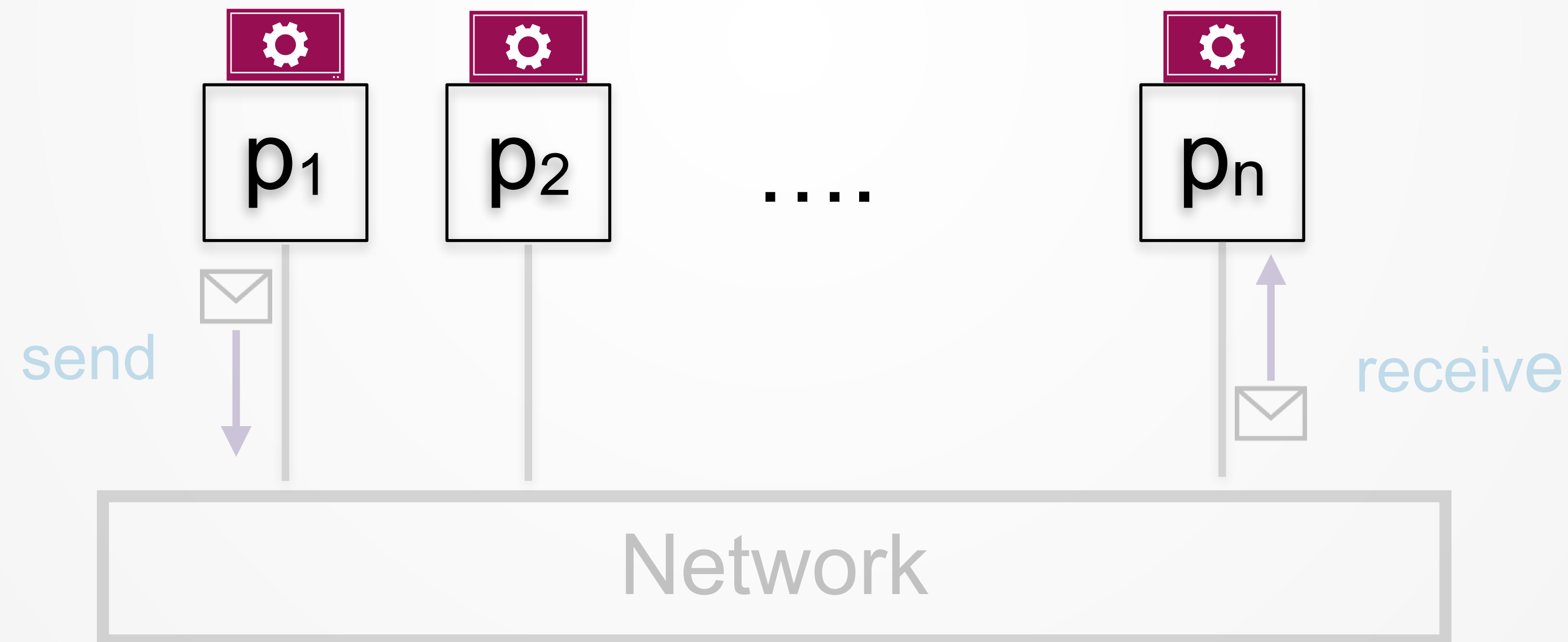
ID2203

KTH-2022

# What is a distributed system?

# WHAT IS A DISTRIBUTED SYSTEM?

"A set of **nodes**, connected by a **network**, which appear to its users as a **single** coherent system"

KTH-2022

# WHAT IS A DISTRIBUTED ALGORITHM

**"A copy of a program running in each process"**

# OUR FOCUS IN THIS COURSE

- Concepts (Processes, Messages, Failures)

- Models (assumptions about system)

- Given the model…
  ‣ Which problems are solvable / not solvable
  ‣ What are the core problems in distributed systems
  ‣ What are the algorithms
  ‣ How to reason about correctness

ID2203

KTH-2022

# WHY STUDY DISTRIBUTED SYSTEMS?

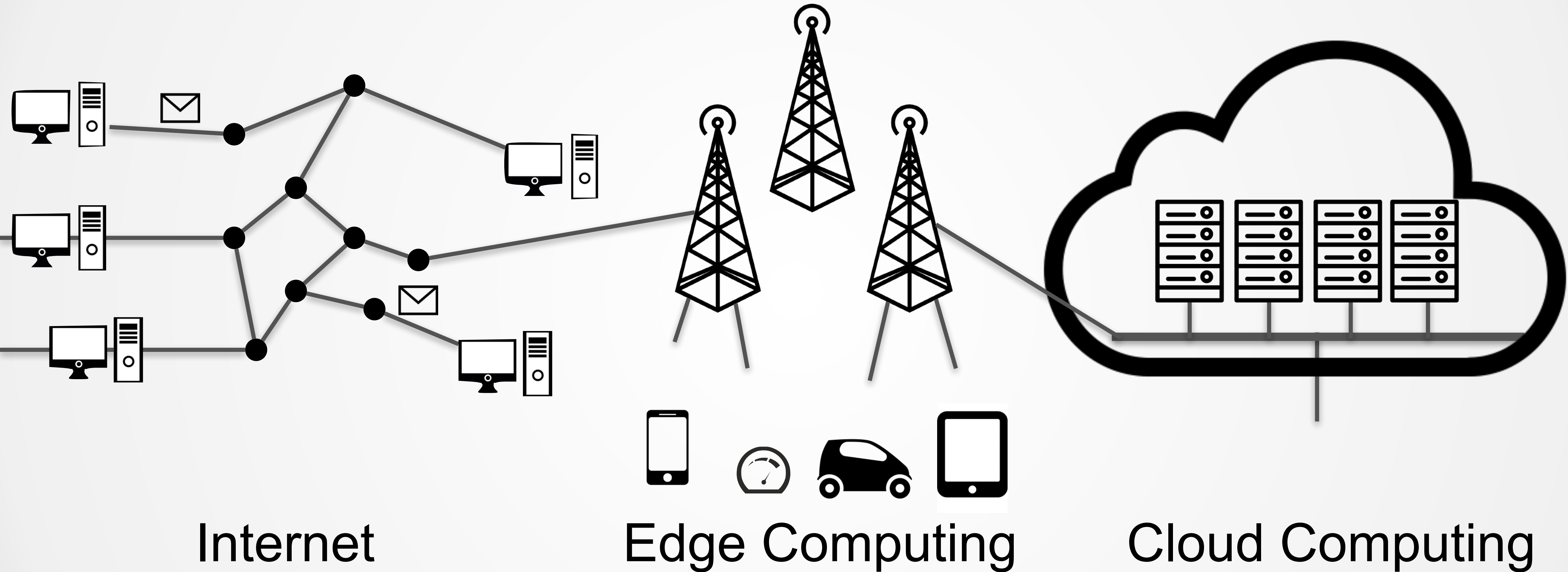It is important, useful and interesting

Societal importance

Internet

WWW

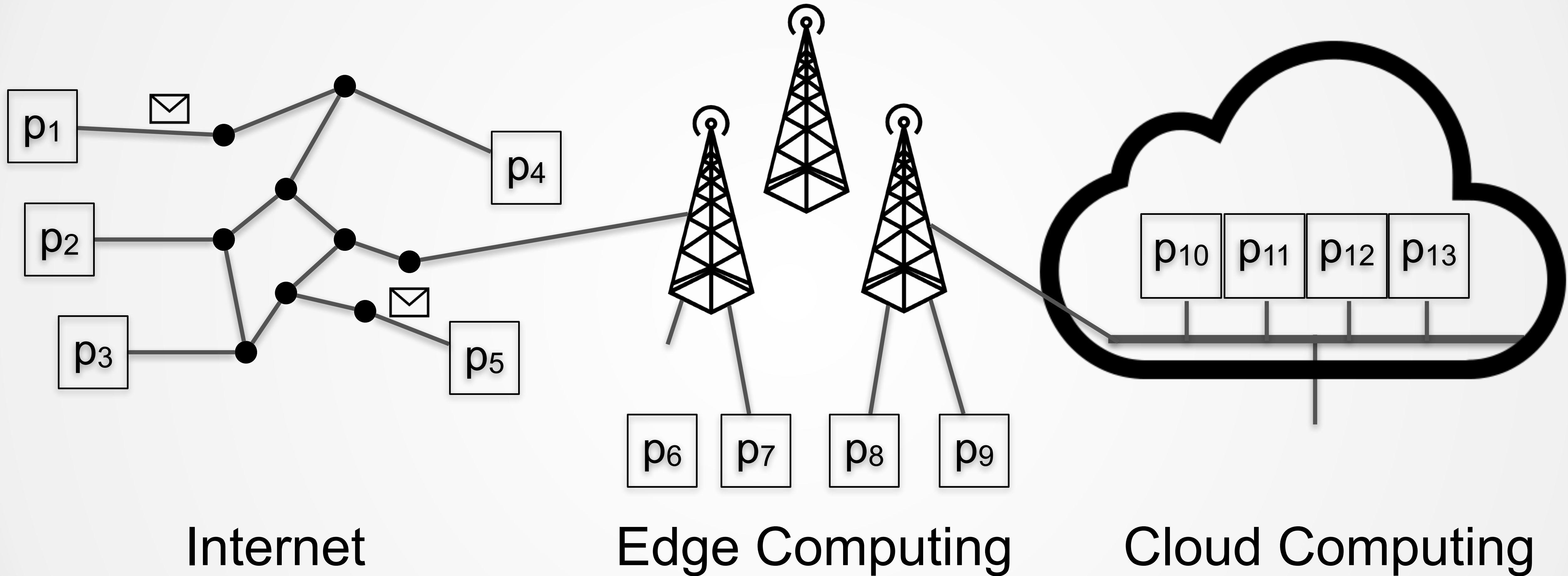Cloud computing (e.g., Google, Amazon)

Edge computing
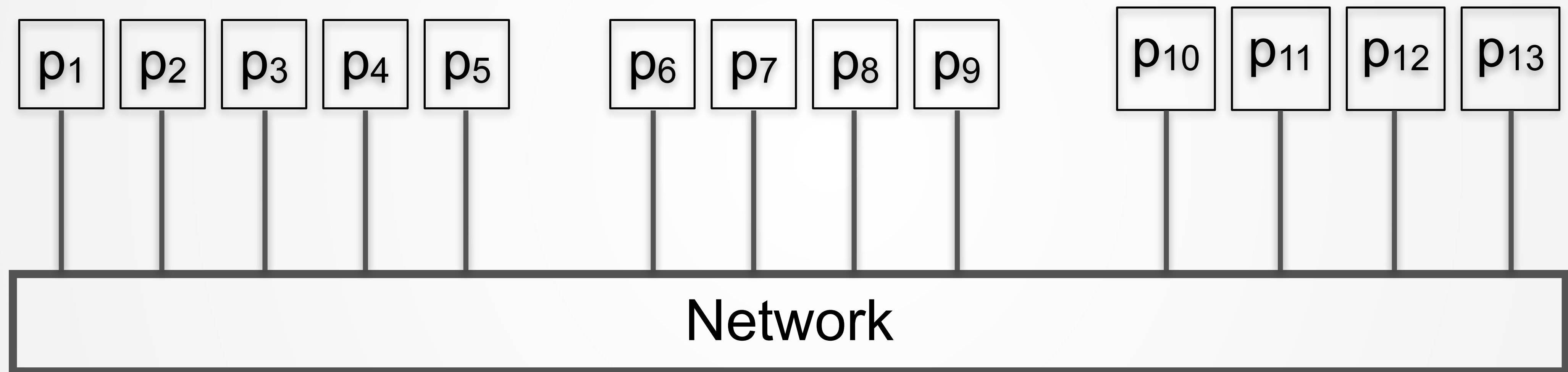
Small devices (mobiles, sensors)

ID2203

KTH

KTH-2022

# WHY STUDY DISTRIBUTED SYSTEMS?



Internet          Edge Computing          Cloud Computing

# WHY STUDY DISTRIBUTED SYSTEMS?



Internet

Edge Computing

Cloud Computing

ID2203

KTH-2022

# WHY STUDY DISTRIBUTED SYSTEMS?

# WHY STUDY DISTRIBUTED SYSTEMS?

It is important and useful

- Technical importance
  - Improve scalability
  - Improve reliability
  - Inherent distribution

ID2203

KTH–2022

# WHY STUDY DISTRIBUTED SYSTEMS?

It is very challenging!

**Partial Failures**

Network (dropped messages, partitions)

Node failures

**Concurrency**

Nodes execute in parallel

Messages travel asynchronously

Parallel computing

Recurring core problems

# Core Problems in Distributed Systems

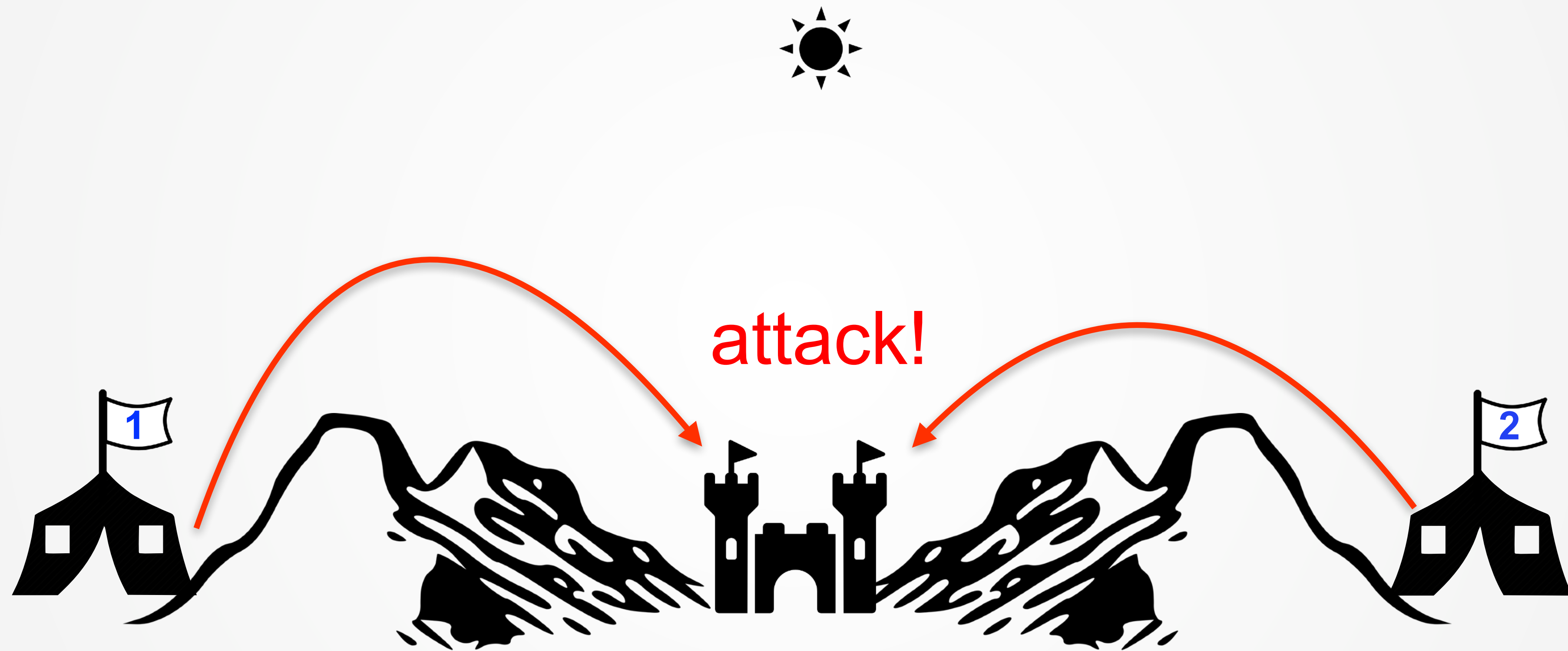## What types of problems are there?

# Teaser: Two Generals' Problem

**"Two generals need to coordinate an attack"**

- Must agree on time to attack
- They'll win only if they attack simultaneously
- Communicate through messengers
- Messengers may be killed on their way

ID2203

KTH-2022

ID2203

KTH–2022

# TEASER: TWO GENERALS' PROBLEM

attack!

1

2

ID2203

KTH-2022

KTH–2022

ID2203

KTH–2022

ambushed!

1

2

Impossible to solve!

ID2203

KTH–2022

# Teaser: Two Generals' Problem

**Applicability to distributed systems**

- Two processes need to agree on a value before a specific time-bound
- Communicate by messages using an unreliable channel

**Agreement is a core problem...**

# CONSENSUS: AGREEING ON A NUMBER

**Consensus problem**
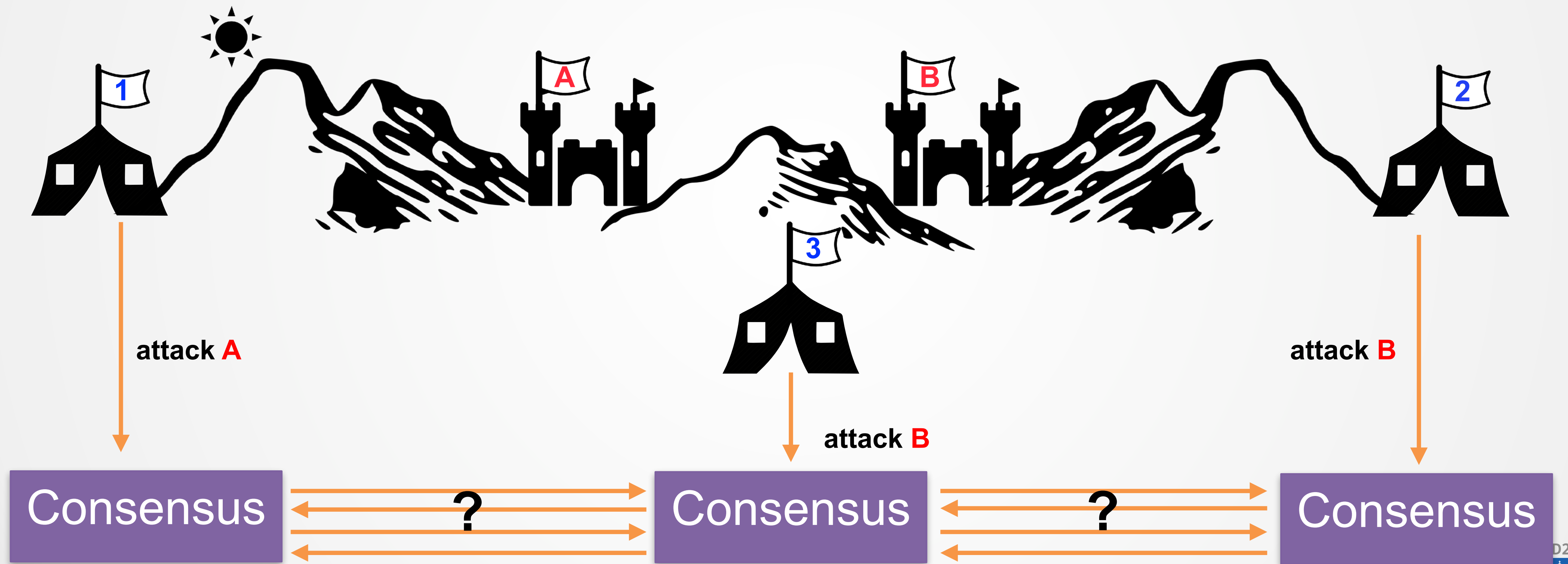
All nodes/processes propose a value
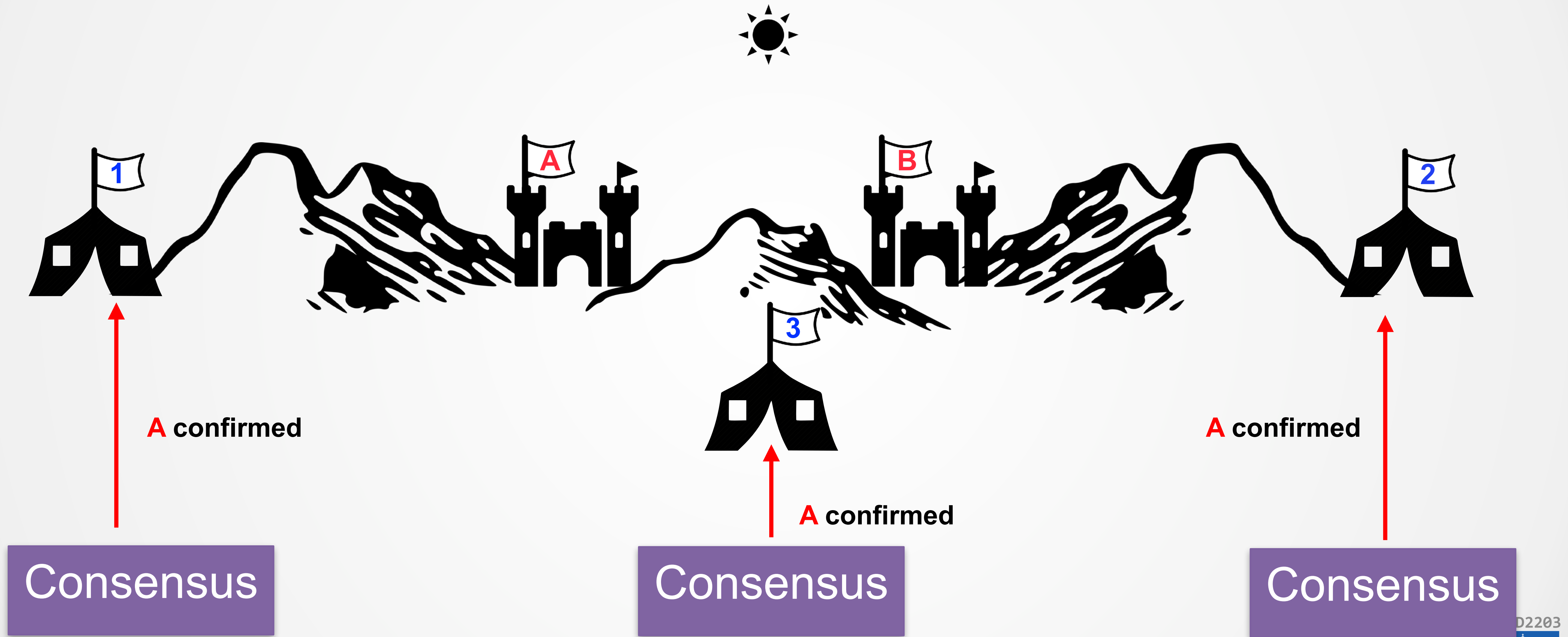Some nodes (non correct nodes) might crash & stop responding

The algorithm must ensure a set of properties (specification):
- ▶ All correct nodes eventually decide
- ▶ Every node decides the same
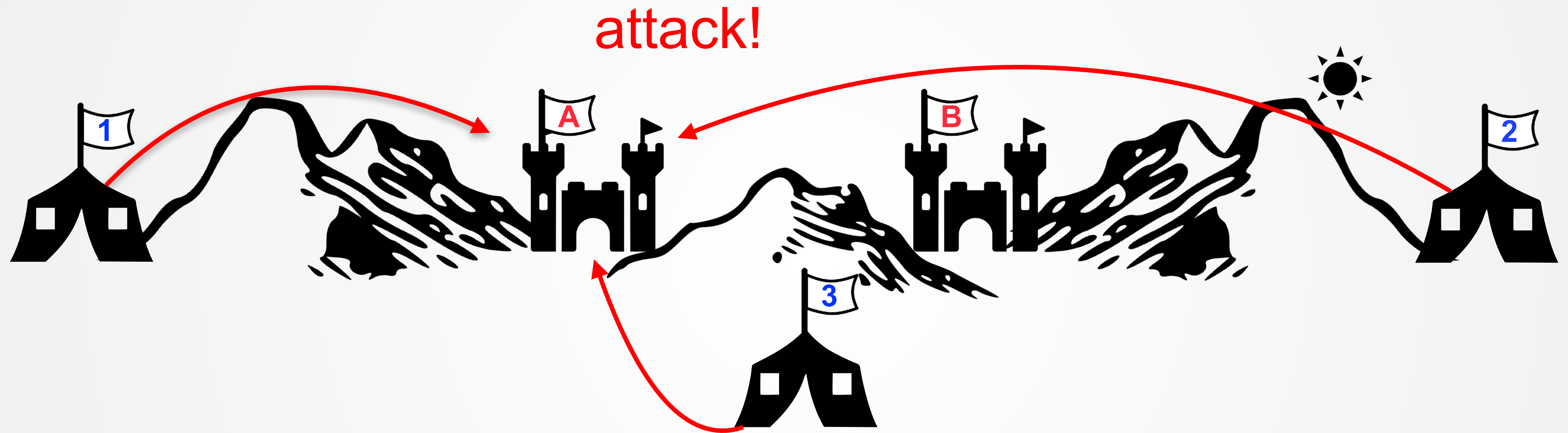- ▶ Only decide on proposed values

ID2203

KTH-2022

# EXAMPLE: AGREEING ON A TARGET

attack!

Consensus

Consensus

Consensus

# EXAMPLE: AGREEING ON A TARGET



ambush!

Consensus ⟷ ? ⟷ Consensus ⟷ ? ⟷ Consensus

# EXAMPLE: AGREEING ON A TARGET

attack!

1

A

B

2

3

Consensus

Consensus

Consensus

D2203

KTH-2022

# Is Consensus Solvable?

Consensus problem
    All nodes propose a value
    Some nodes might crash & stop responding

The algorithm must ensure:
▶ All correct nodes eventually decide
▶ Every node decides the same
▶ Only decide on proposed values

ID2203

KTH–2022

**Distributed Databases / Cloud Stores**

Concurrent changes/transactions to same data

Nodes should <span style="color:#8a1538">agree</span> on changes

Use a kind of consensus: <span style="color:#8a1538">atomic commit</span>

Only two proposal values <span style="color:#1f6fa8">{commit, abort}</span>

# BROADCAST PROBLEM

## Atomic Broadcast

▸ A node broadcasts a message

▸ If sender correct, all correct nodes deliver msg

▸ All correct nodes deliver the same messages

▸ Messages delivered in the same order

# ATOMIC BROADCAST IS IMPORTANT

**Replicated services**

- Multiple servers (processes)
- Execute the same sequence of commands
- Replicated State Machines RSM

Use atomic broadcast
Provide fault tolerance

# Can we use atomic broadcast to <u>solve</u> consensus?

# ATOMIC BROADCAST ⟷ CONSENSUS

I. **Atomic broadcast** can be used to solve **Consensus**!
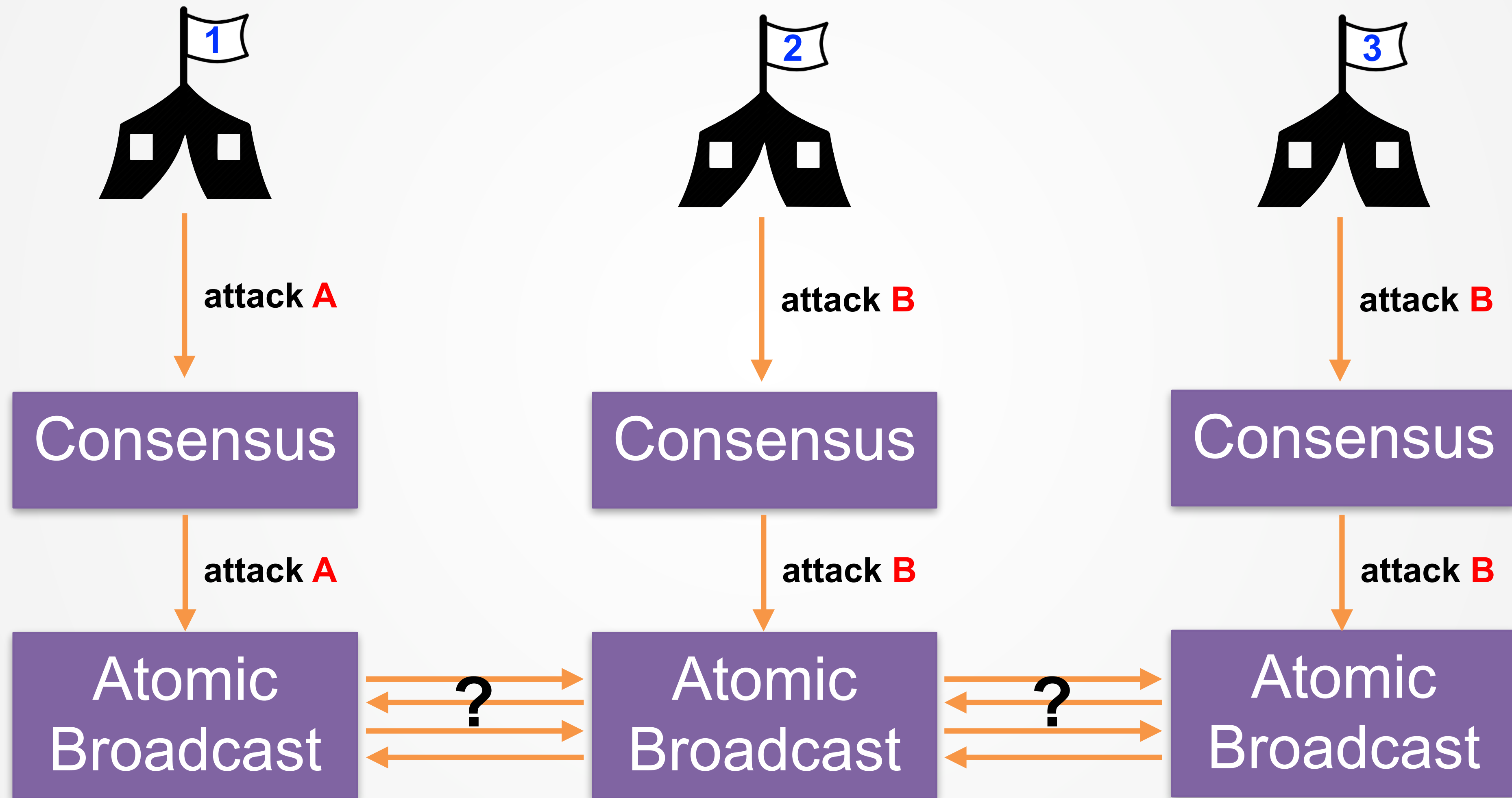

i.e., Every node broadcasts its proposal
- Decide on the **first** received proposal

- Messages received in same order
    - Thus, all nodes will decide the same value.


II. **Consensus** can be used to solve **Atomic broadcast**
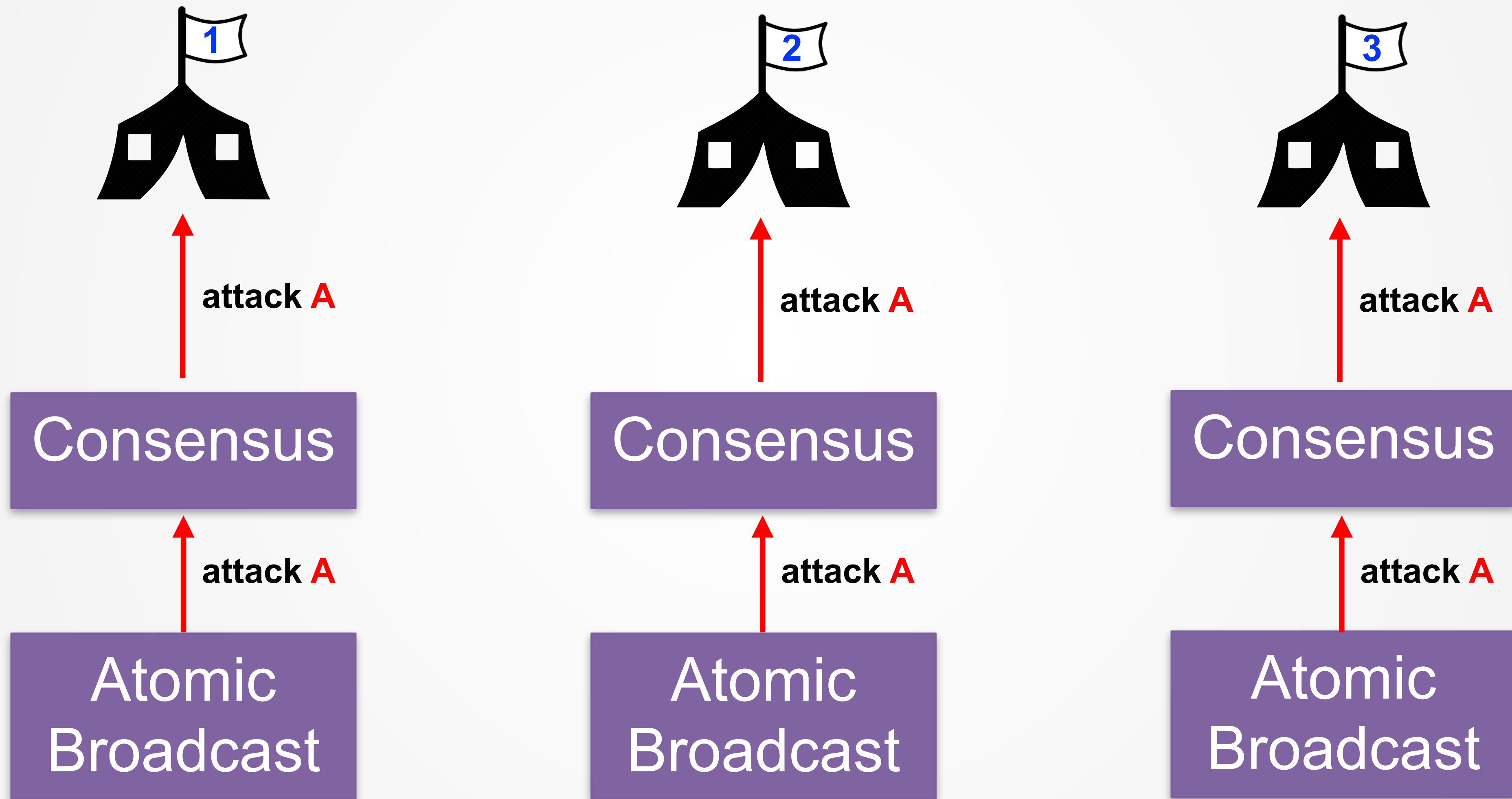
**(more on that later in the course)**


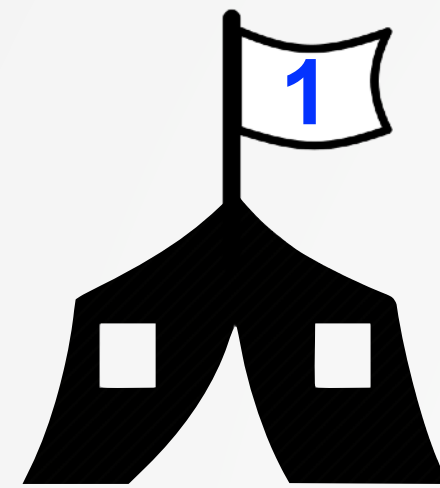I+II: Atomic Broadcast **equivalent** to Consensus

# ATOMIC BROADCAST ↔ CONSENSUS



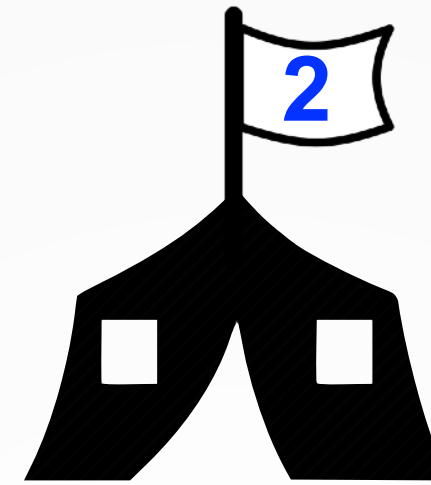attacking **A**    attacking **A**    attacking **A**

Consensus    Consensus    Consensus

↑ attack **B**    ↑ attack **B**    ↑ attack **B**

Atomic Broadcast    Atomic Broadcast    Atomic Broadcast

ID2203
KTH-2022

# **Models** of Distributed Systems

## How to reason about them?

# MODELLING A DISTRIBUTED SYSTEM

**Timing assumptions**

**Processes**

▸Bounds on time to make a computation step

**Network**

▸Bounds on time to transmit a message between a sender and a receiver

**Clocks**

▸Lower and upper bounds on clock drift rate

ID2203

KTH-2022

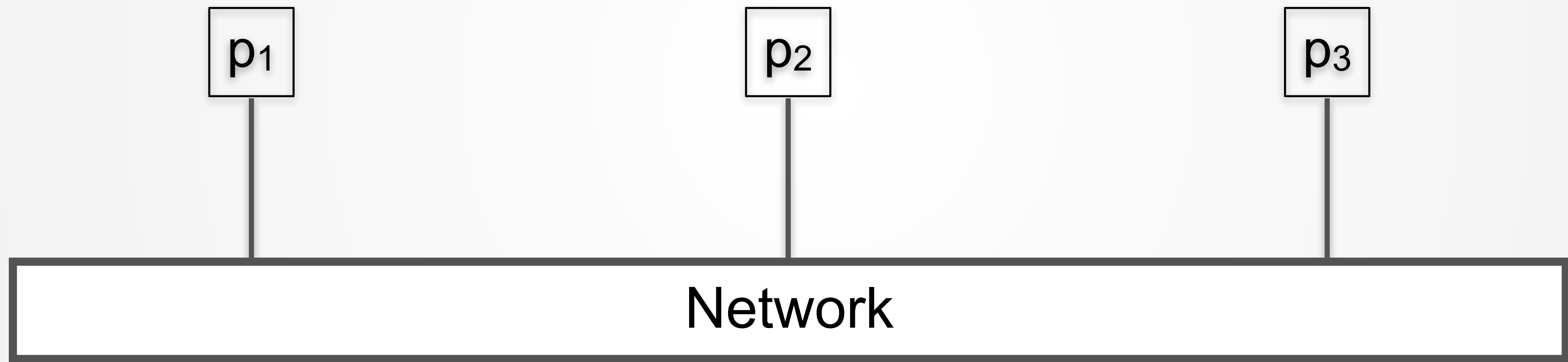# Modelling a Distributed System

**Failure assumptions**

**Processes**

▸ What kind of failure a process can exhibit?

▸ Crashes and stops

▸ Behaves arbitrary (Byzantine)

**Network**

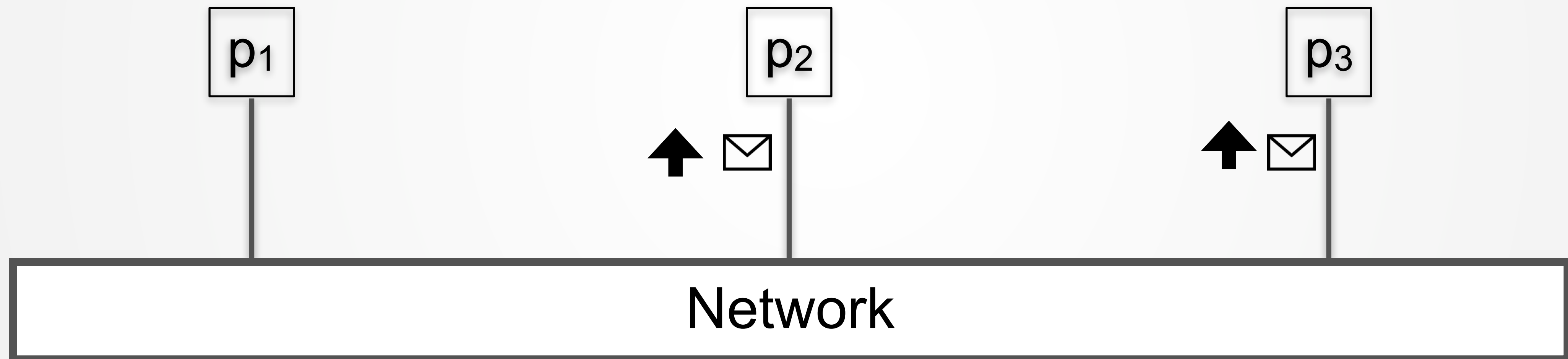▸ Can a network channel drop messages?
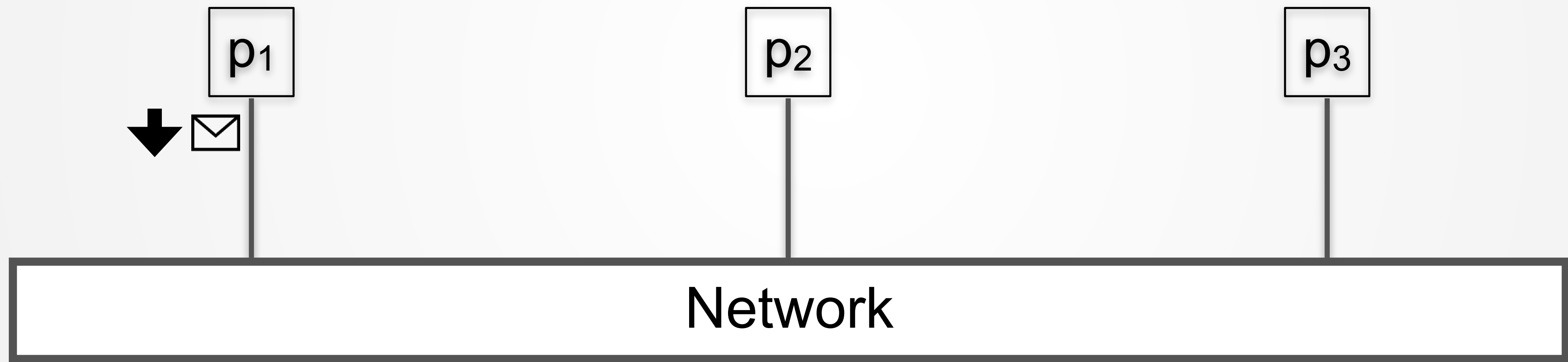
ID2203

KTH-2022

# MODELING A DISTRIBUTED SYSTEM

# Modeling a Distributed System

# Modeling a Distributed System

Network

# NETWORK FAILURES

# PROCESS FAILURES

# PROCESS FAILURES

# PROCESS FAILURES



p₁

p₂

p₃

Network

# PROCESS FAILURES

# BYZANTINE PROCESSES

# BYZANTINE PROCESSES



p₁     p₂     p₃

Network

# BYZANTINE PROCESSES

# MODELLING A DISTRIBUTED SYSTEM

**The Asynchronous System Model**

▸No bound on time to deliver a message

▸No bound on time to compute

▸Clocks are not synchronized

Internet is essentially asynchronous

ID2203

KTH-2022

# IMPOSSIBILITY OF CONSENSUS

**Consensus** <u>**cannot**</u> be solved in asynchronous system if node crashes can happen.

Implications on

▸Atomic broadcast

▸Atomic commit

▸Leader election

...

ID2203

KTH-2022

# MODELLING A DISTRIBUTED SYSTEM

**Synchronous system**

▸Known bound on time to deliver a message (latency)

▸Known bound on time to compute

▸Known lower and upper bounds in physical clock drift rate

Examples:

▸Embedded systems (shared clock)

▸Multicore computers

ID2203

KTH–2022

# POSSIBILITY OF CONSENSUS

**Consensus** is **solvable** in synchronous system with up to N-1 crashes

Intuition behind solution

▸Accurate crash detection

  ▸Every node sends a message to every other node

  ▸If no msg from a node within bound, node has crashed

Not useful for Internet, how to proceed?

ID2203

KTH-2022

# MODELLING A DISTRIBUTED SYSTEM

A more realistic view of most systems (e.g., over internet)

▸ Bounds respected mostly

▸ Occasionally violate bounds (congestion/failures)

How do we model this?

**Partially synchronous system**

▸ Initially system is asynchronous

▸ Eventually the system becomes synchronous

Consensus **solvable** in partially synchronous system

with up to N/2 crashes

.Can't this be used in **Cloud services**?

KTH-2022

# FAILURE DETECTORS

Let each node use a **failure detector**

▸Detects crashes

▸Implemented by heartbeats and waiting

▸Might be initially wrong, but eventually correct

Consensus and Atomic Broadcast solvable with failure detectors

How? Attend rest of course!

ID2203

KTH-2022

# MODELING A DISTRIBUTED SYSTEM

**Timed Asynchronous system**

▸ No bound on time to deliver a message

▸ No bound on time to compute

▸ Clocks have known clock-drift rate

Another realistic model for the Internet

ID2203

KTH-2022

# BYZANTINE FAULTS

Some processes might behave arbitrarily
- ▶ Sending wrong information
- ▶ Omit messages…

Byzantine algorithms that tolerate such faults
- ▶ Only tolerate up to **1/3** Byzantine processes
- ▶ Non-Byzantine algorithms can often tolerate **½** nodes in the asynchronous model

ID2203

KTH-2022

# SELF-STABILIZING ALGORITHMS

**Wont be covered in the course but cool to know.**

▸Robust algorithms that run forever
System might temporarily be incorrect
But eventually always becomes correct

▸ System can either be in a legitimate state or an illegitimate
state (invariant)

Self-stabilizing algorithm iff
Convergence
Given any illegitimate state, system eventually goes to a legitimate state
Closure
If system in a legitimate state, it remains in a legitimate state

ID2203

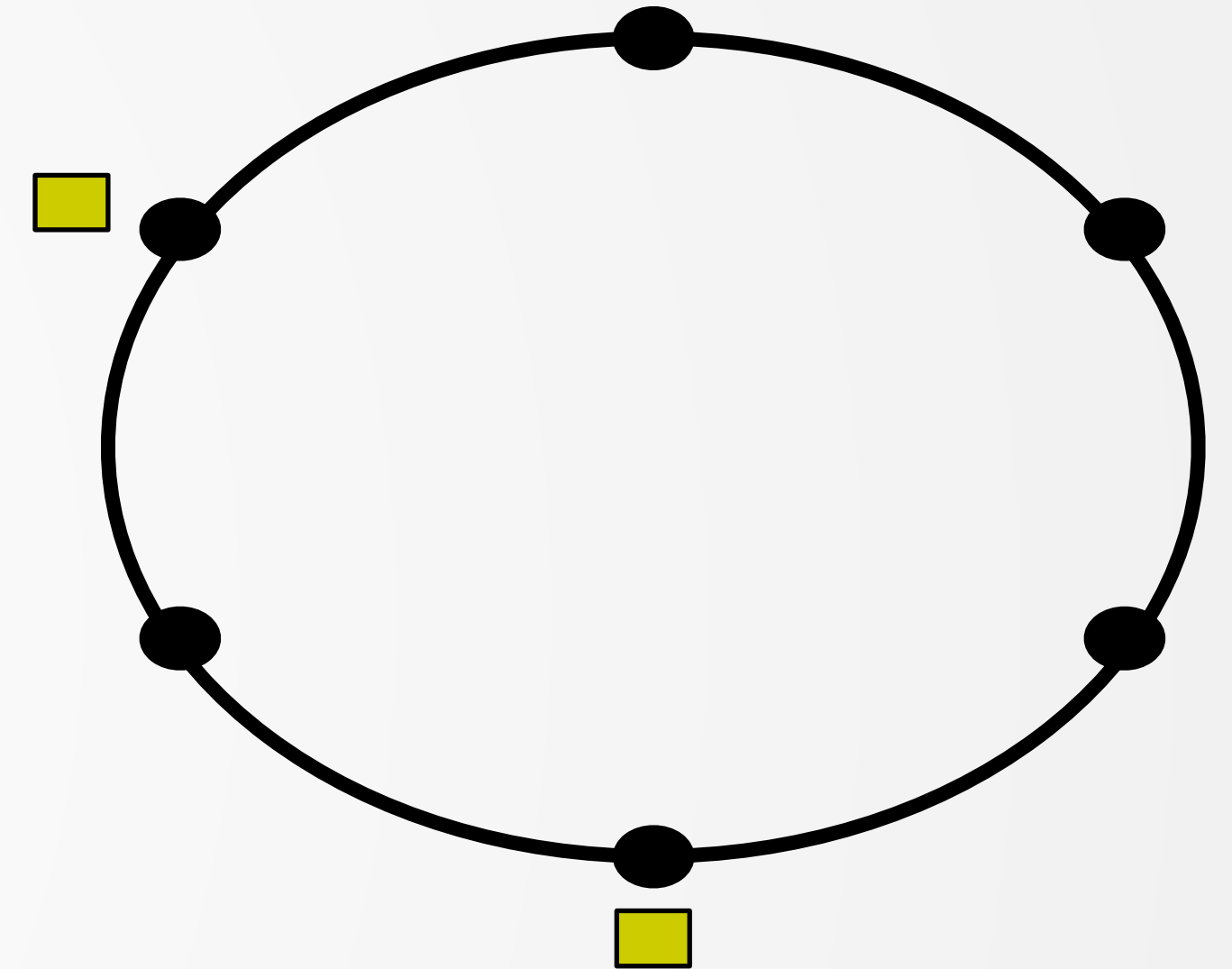KTH-2022

# SELF-STABILIZING EXAMPLE

## Token ring algorithm

Wish to have one token at all
times circulating among processes

## Self-Stabilization

Error leads to 2,3,… tokens
Ensure always 1 token eventually

S. Haridi, KTHx ID2203x

ID2203

KTH–2022

# SUMMARY

**Distributed systems everywhere**

Set of processes (nodes) cooperating over a network

Few **core problems** reoccur

Consensus, Broadcast, Leader election, Shared Memory

**Different failure scenarios** important

Crash stop, Byzantine, self-stabilizing algorithms

Interesting **research** directions

Large scale dynamic distributed systems

ID2203

KTH-2022