

# Fault-Tolerant Distributed Systems

Rui Carlos Oliveira

[rco@di.uminho.pt](mailto:rco@di.uminho.pt)

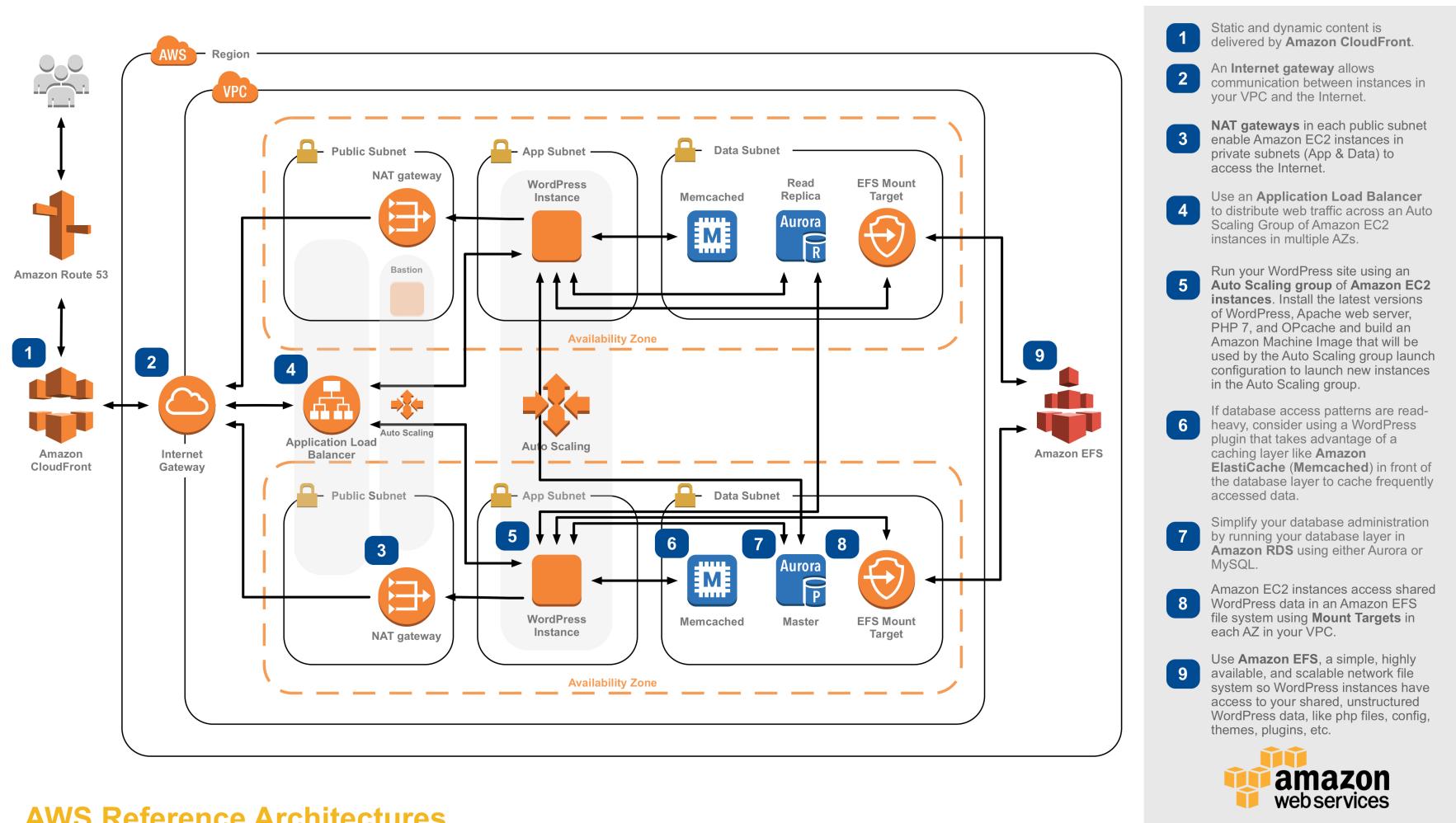
Informatics Department  
Minho University

2022/23

- **Introduction to fault-tolerant distributed systems**
- Models of distributed systems and related faults
- Data replication
- Distributed consensus
- State machine replication
- Database replication

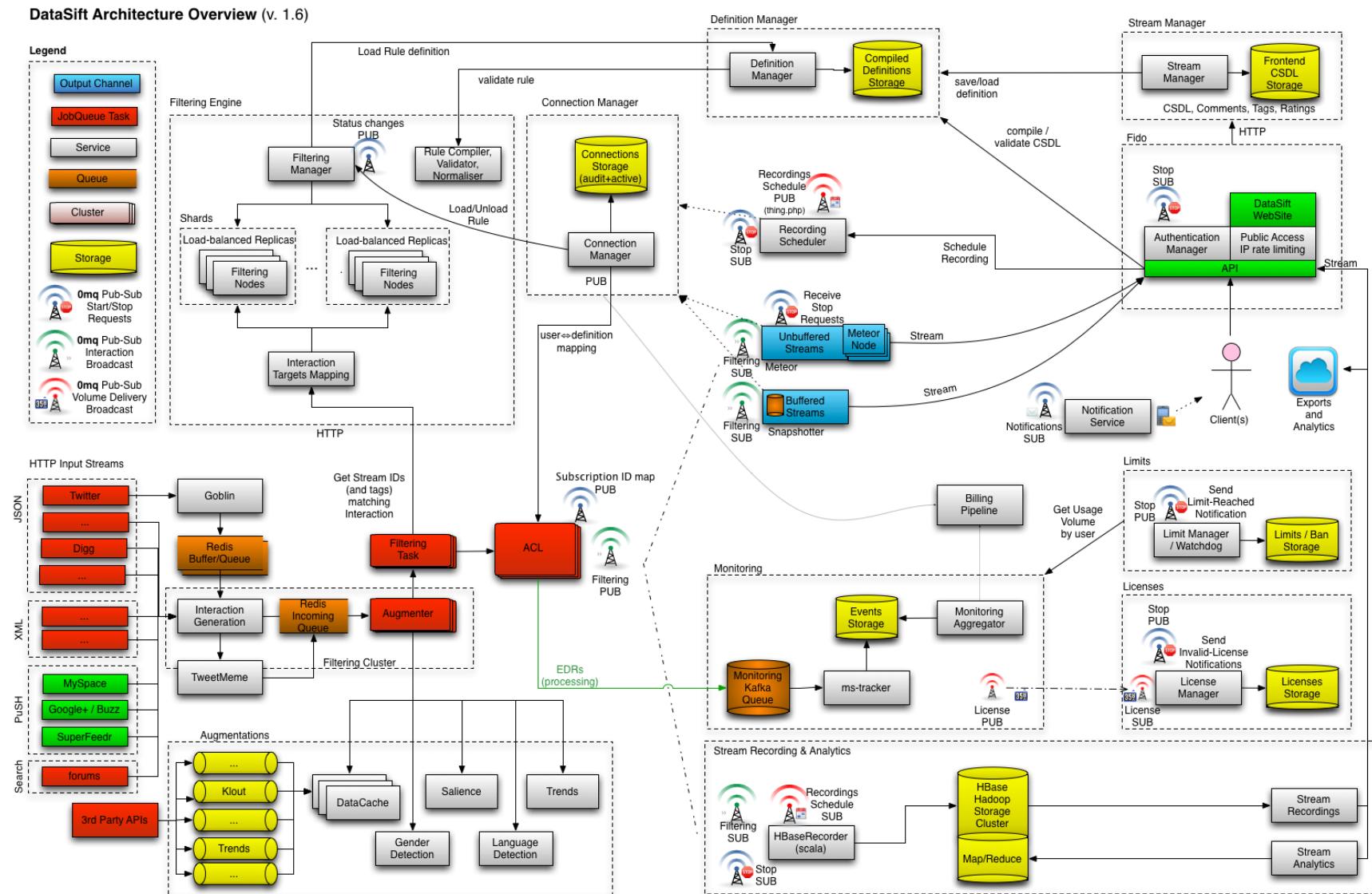
# Introduction to fault-tolerant distributed systems

- Software systems are ever more complex, bigger, and full of intricate dependencies



# Introduction to fault-tolerant distributed systems

- Software systems are ever more complex, bigger, and full of intricate dependencies



# Introduction to fault-tolerant distributed systems

- ▶ Goal of the course

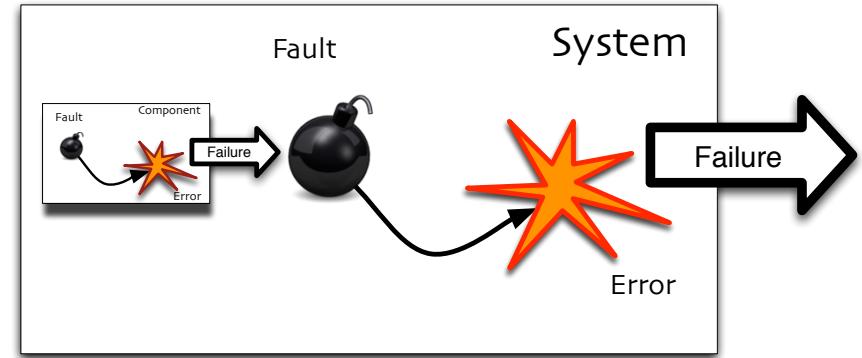
- ▶ Learn how to make a software system fault-tolerant

- ▶ Impact of faults

- ▶ Performance
  - ▶ Unavailability
  - ▶ Misbehaviour

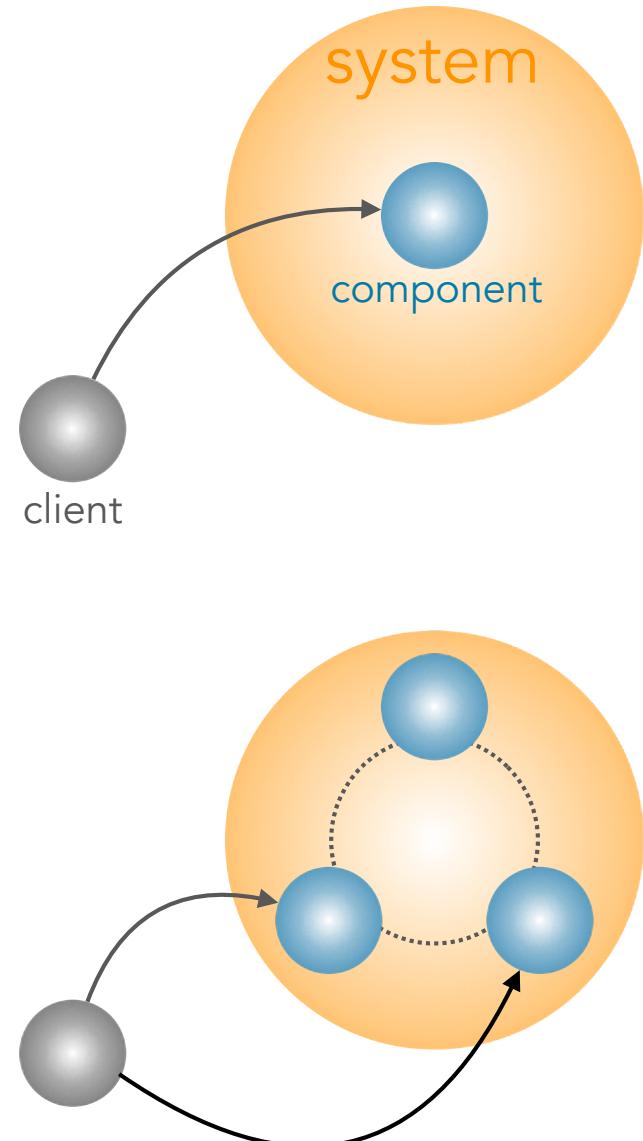
- ▶ The challenge

- ▶ Make *fault-tolerant X* service as available as possible (availability)
    - ▶ 99% to 99.999% of uptime (3.65 days to 5.2 minutes of downtime)
  - ▶ Any behaviour of *fault-tolerant X*, needs to be a behaviour of X (safety)



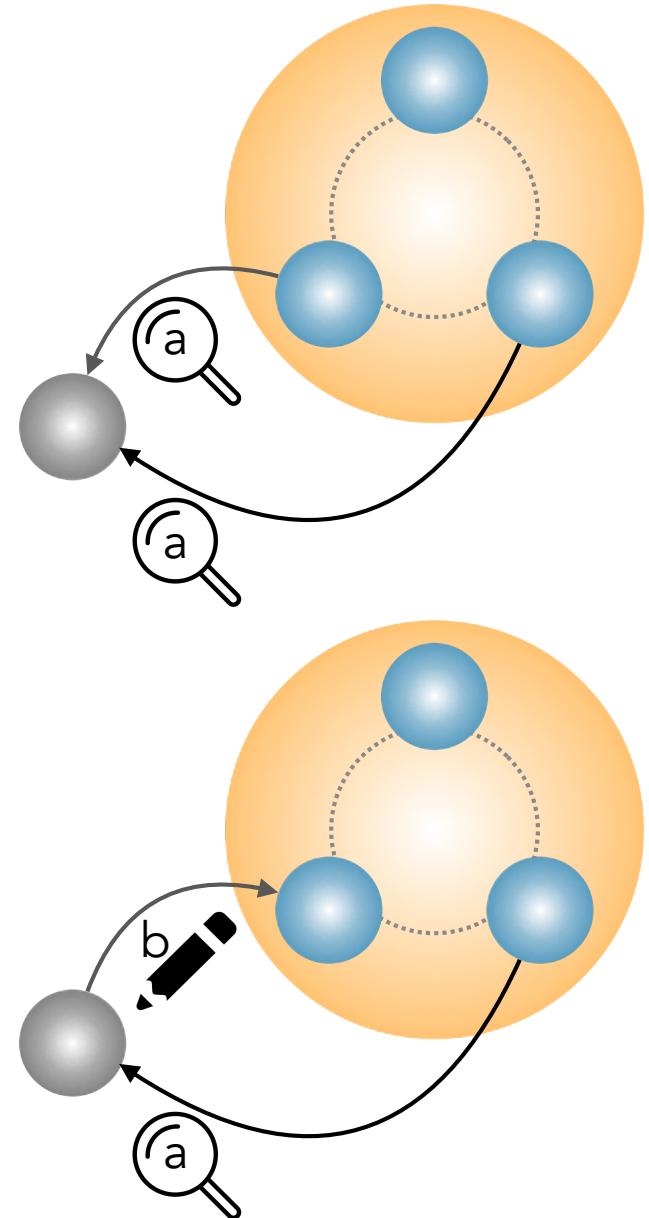
# Introduction to fault-tolerant distributed systems

- Fault tolerance requires replication
  - As any system component is prone to fail
  - Any vital component needs to be replicated
    - Remove any *Single Point Of Failure (SPOF)*
- Distributed systems suit replication well
  - Modularity
  - Incremental redundancy
  - Graceful degradation
  - Heterogeneity
  - Encapsulation



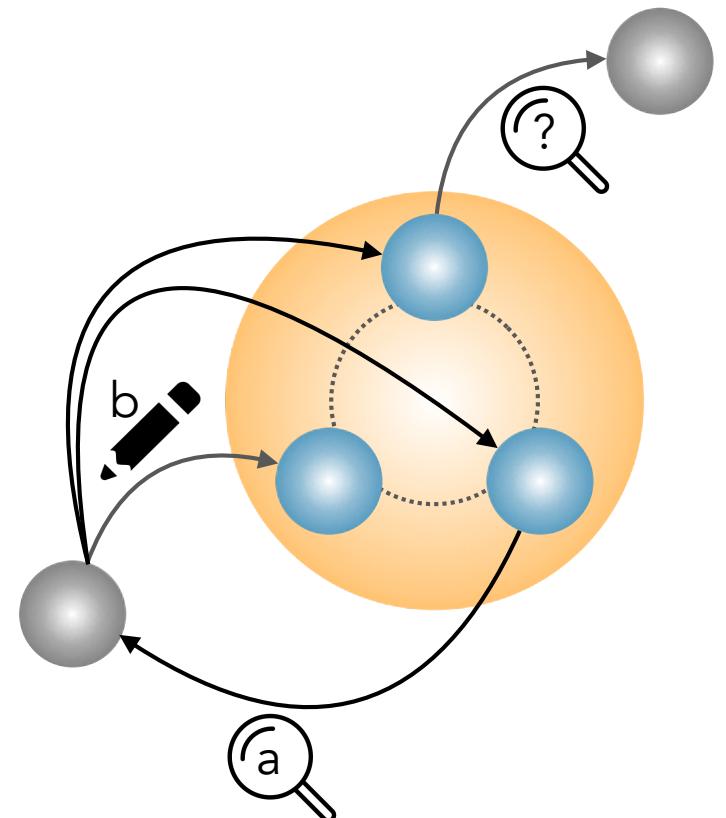
# Introduction to fault-tolerant distributed systems

- Any replication approach depends on
  - The component's state mutability



# Introduction to fault-tolerant distributed systems

- Any replication approach depends on
  - The component's state mutability
  - The concurrency model



# Introduction to fault-tolerant distributed systems

- Any replication approach depends on
  - The component's state mutability
  - The concurrency model
  - The type and number of faults (fault model)
  - The workload bias (usually between reads and writes)

# Introduction to fault-tolerant distributed systems

## Reading material

- A. Avizienis, J-C Laprie, B. Randell, and C. Landwehr

"Basic Concepts and Taxonomy of Dependable and Secure Computing",  
IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING, VOL. 1, NO. 1, JANUARY-MARCH 2004



- Introduction to fault-tolerant distributed systems
- **Models of distributed systems and related faults**
- Data replication
- Distributed consensus
- State machine replication
- Database replication

# Models of distributed systems and related faults

## • Experimental observation

- Build things and observe how they behave in various situations
- Experience enables us to build things for situations that have not been studied

Practice

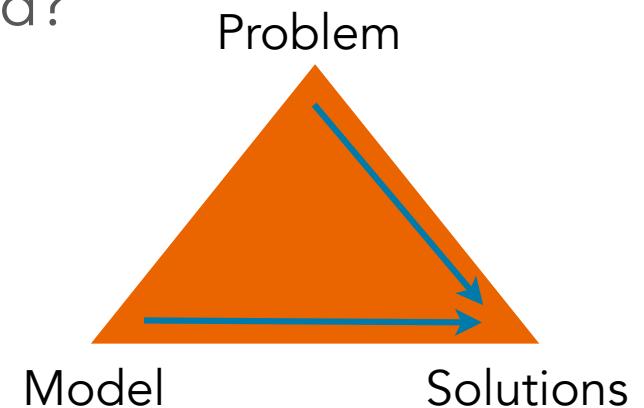
## • Modelling and analysis

- Formulate a model by simplifying the object of study and defining a set of rules to define its behaviour
- Analyze the model and infer consequences

Theory

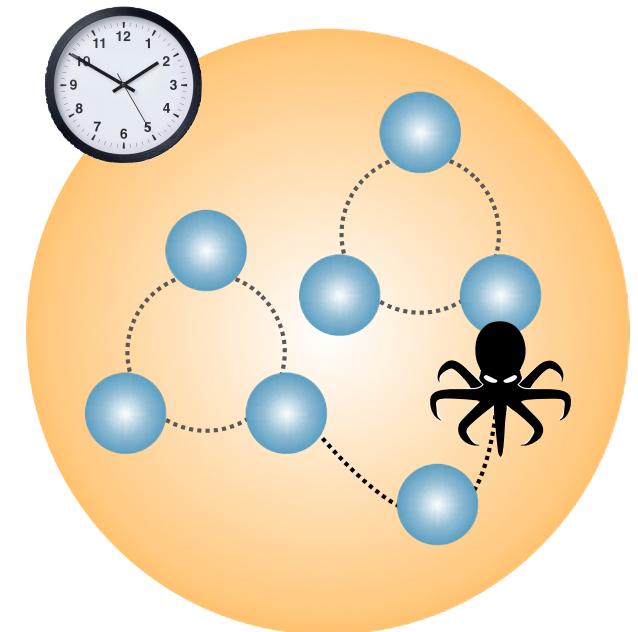
# Models of distributed systems and related faults

- Modelling a system consists of identifying and precisely characterising its entities and the behavioural rules that are relevant for the problems in hand
- A model needs to be tractable in that it does not contain uninteresting details that would impair the desired analysis and, at the same time, it needs to be accurate so that it captures all attributes affecting the phenomena of interest
- We seek answers to two fundamental questions:
  - Feasibility: What kind of problems can be solved?
  - Cost: What is the relative cost of solutions?



# Models of distributed systems and related faults

- A minimal model for a distributed system consists of a finite set of autonomous sequential processes connected through a finite set of point-to-point communication channels
- We picture the existence of an adversary with a limited power to impair the correct behaviour of processes or channels
- For the problems we will be considering we model the system as:
  - A finite set of sequential processes
  - A finite set of communication point-to-point channels
  - An adversary
  - Rules governing the behaviour of the above



# Models - What kind of problems can be solved?

- System model:

- Two processes, A and B, communicate by sending and receiving messages
- Neither process can fail. However, the channel can experience transient failures, resulting in the loss of a subset of the messages that have been sent



- A **Coordination problem**:

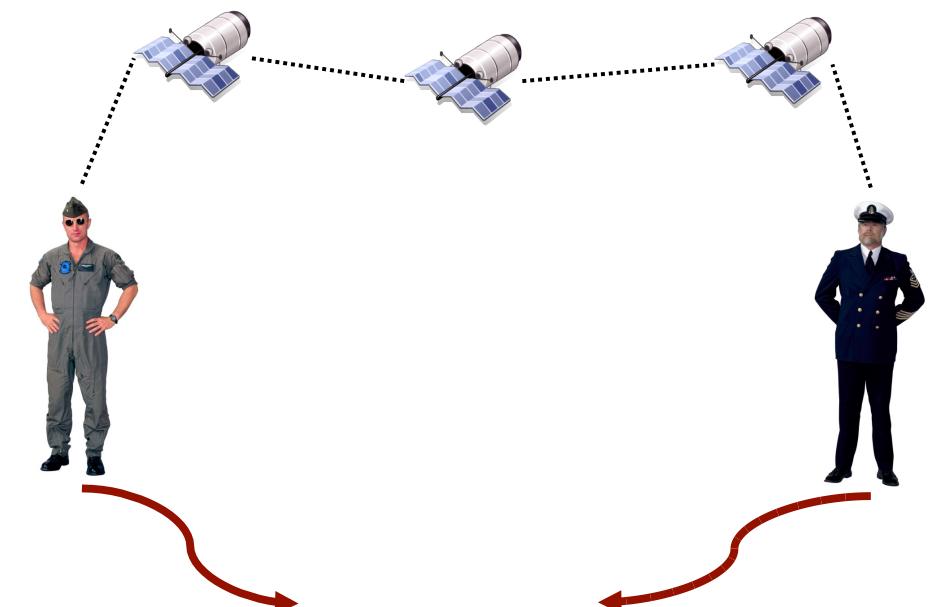
- Devise a protocol where either of two actions  $\alpha$  and  $\beta$  are possible, but (i) both processes take the same action and (ii) neither takes both actions

Problem: Agree on a simultaneous attack

Both attack: Win!

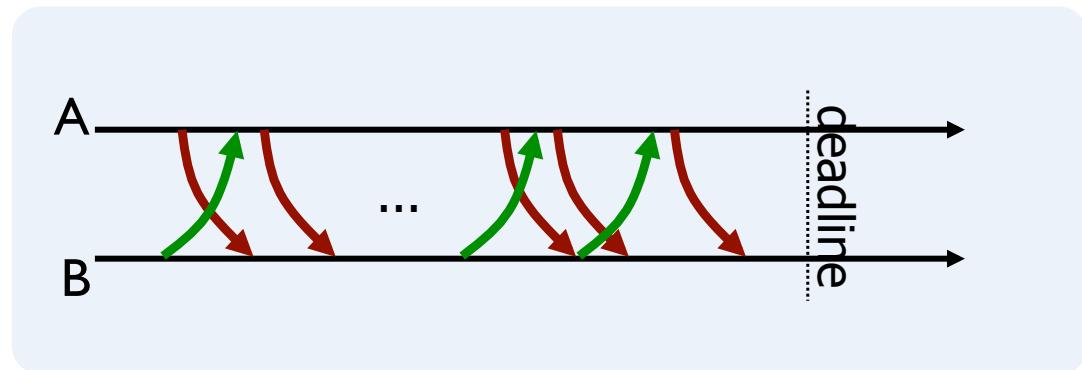
Only one attacks: Dead

No one attacks: Safe but...



# Models - What kind of problems can be solved?

- Any solution will consist in a message exchange between A and B
- Without channel failures the solution would be straightforward
- Interesting case: Messages received is a subset of messages sent
- Consider an hypothetical solution:



- Does the last message matter?
- And the second last?
- What if the “solution” was the one with the minimum message exchange?

# Models - Synchrony and Asynchrony

- Of major importance when modelling a distributed system are the assumptions we make regarding the synchrony or asynchrony of processes and communication channels
- Modelling an attribute as asynchronous means that we make no assumptions, absolute or relative, with respect to the time it takes to perform an action. An asynchronous model leads to universal solutions with respect time
- An attribute is synchronous if we assume it take at most  $t$  time units to perform an action
- Models can be semi-synchronous in that the attributes are synchronous
  - But  $t$  is unknown
  - But only eventually

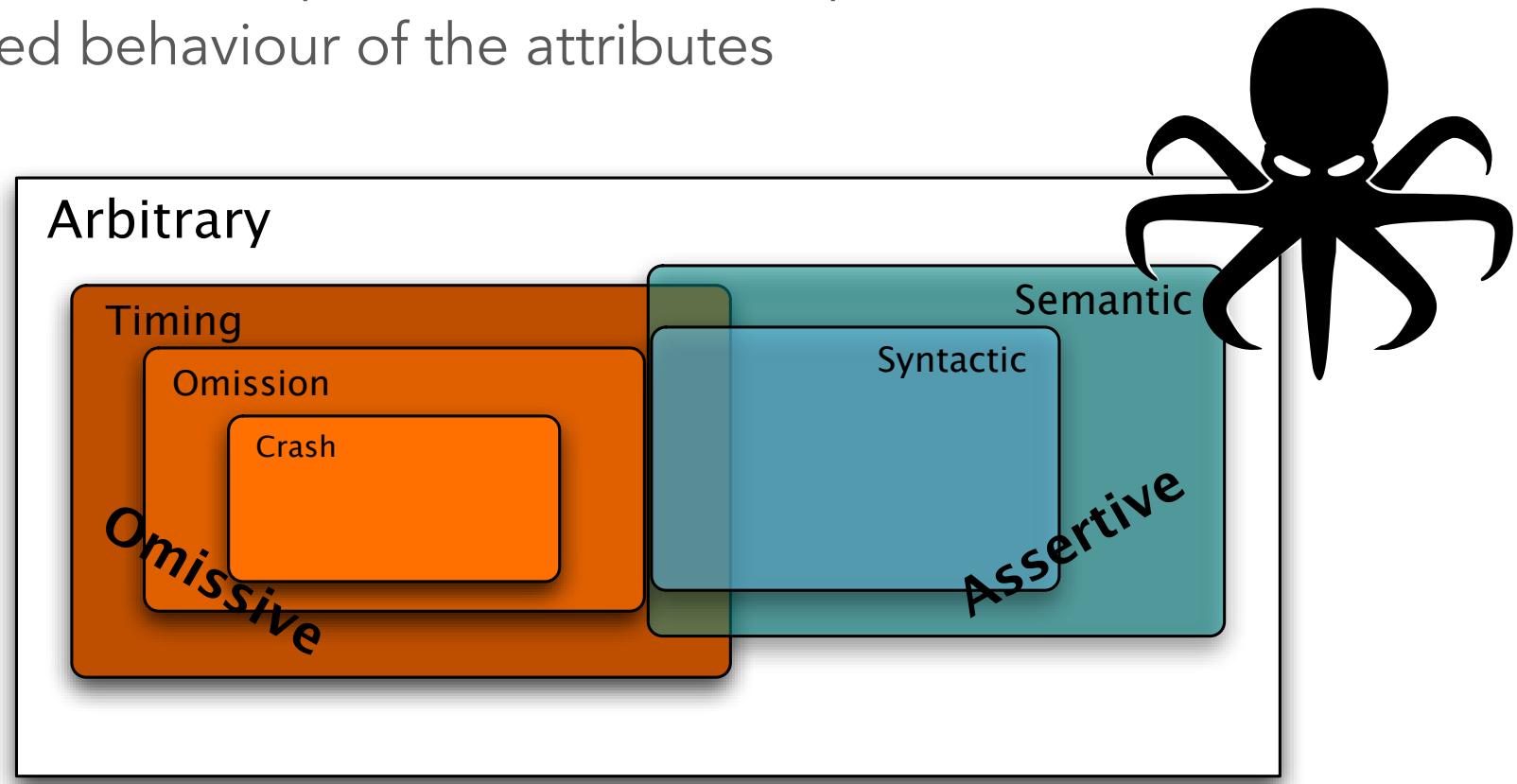


# Models - What is the relative cost of solutions?

- System model A:
  - A set of processes  $P_1, \dots, P_n$ . Each process  $P_i$  has a unique integer identifier  $id_i$
  - There is no adversary
  - All processes start executing at the same time but **no synchrony assumptions**
- An Election problem:
  - Devise a protocol to elect a unique leader eventually known by all processes
- System model S:
  - Processes take at most  $t_p$  per step and messages take at most  $t_m$  to be delivered
- Homework:
  - Devise a solution in S that is not a solution in A. How do they compare?

# Models of distributed systems and related faults

- The characterisation of the adversary can be done through the identification of the type, number and frequency of the deviations to the specified behaviour of the attributes



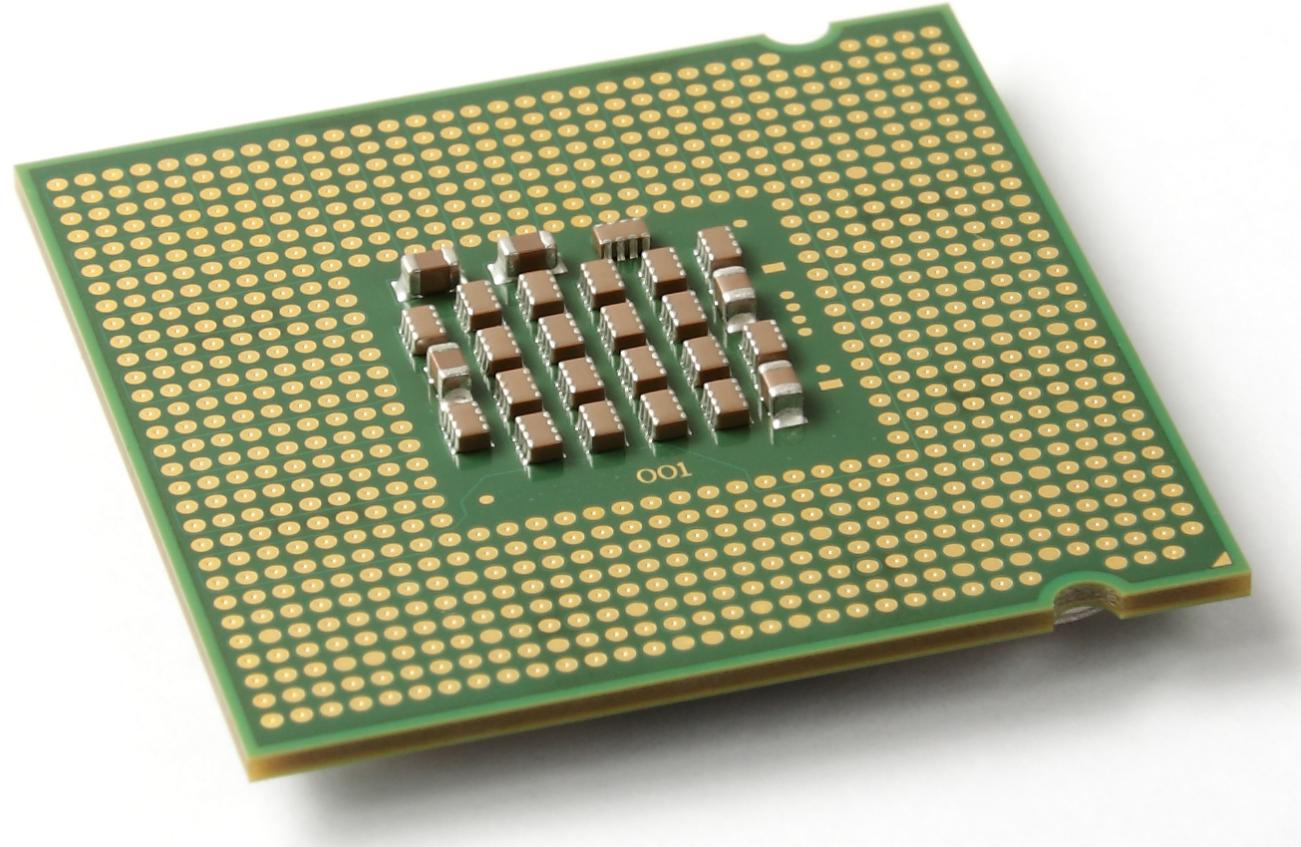
- Omissive: fail-stop, crash-stop, crash-recovery, crash-link, receive, send and general omissions
- Assertive: Syntactic/Semantic, Byzantine

- ▶ Set of processes

- ▶ Unknown
- ▶ Known size
- ▶ Known ids

- ▶ Process faults

- ▶ Crash
- ▶ Crash and recovery
- ▶ Misbehave arbitrarily



# Models of distributed systems - Communication

- Topology:

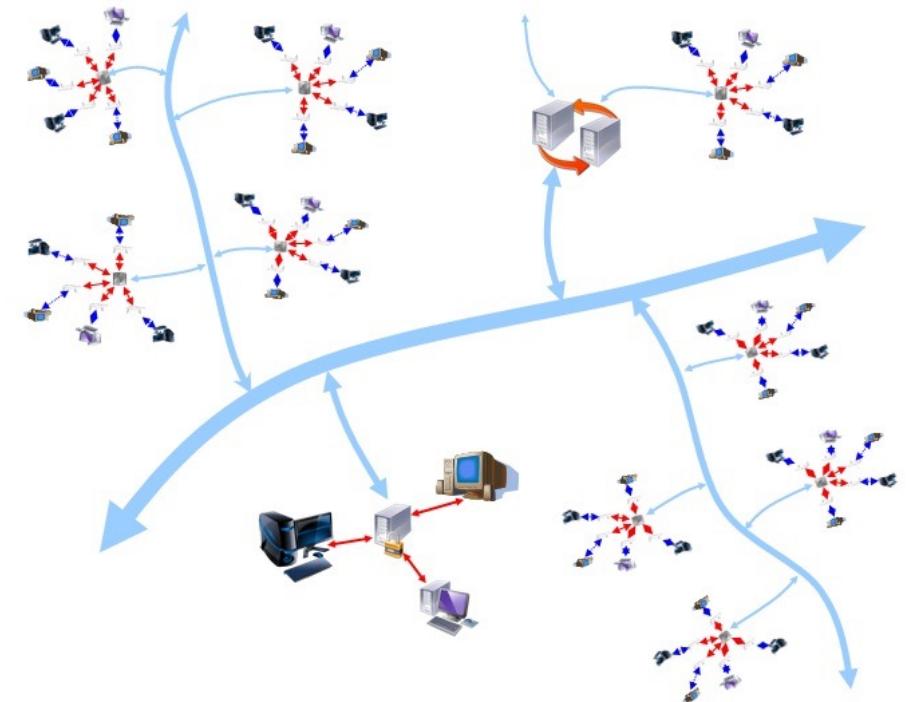
- Fully connected (clique)

- Connected graph:

- Static or Dynamic
  - Partitions
    - Transient or Permanent

- Reliability:

- From all messages lost to no messages lost
  - Can messages be created?



- ▶ Synchronous:

- ▶ Known upper bounds on message delays, processing delays, and clock differences

- ▶ Partially synchronous:

- ▶ Unknown fixed bounds
  - ▶ Known eventual bounds

- ▶ Asynchronous:

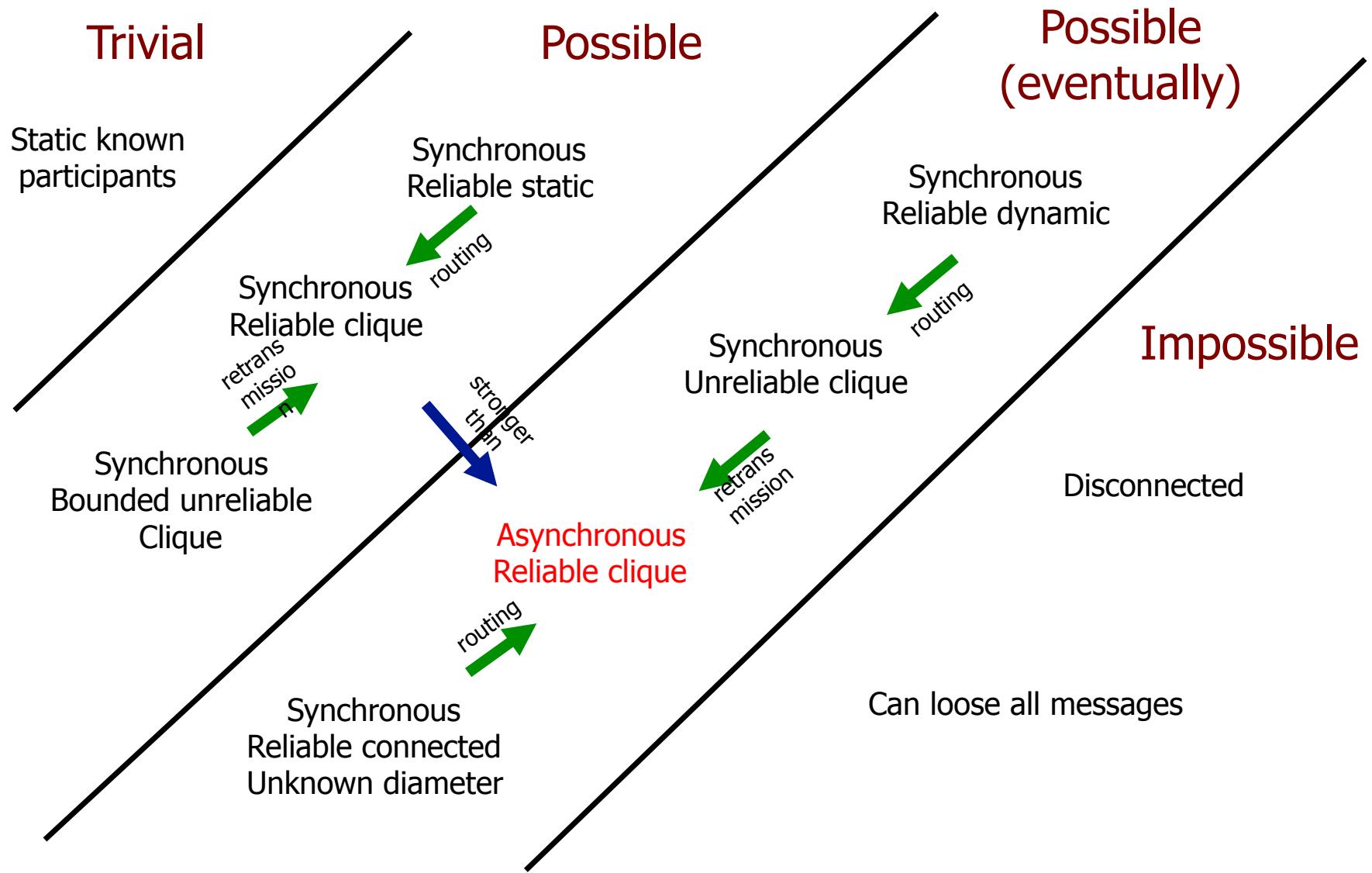
- ▶ No assumptions on bounds and clocks



# Models - What kind of problems can be solved?

|   |   |                                  |
|---|---|----------------------------------|
| Static known participants                   | Synchronous<br>Reliable static                        | Synchronous<br>Reliable dynamic  |
|   | Synchronous<br>Reliable clique                        | Synchronous<br>Unreliable clique |
| Synchronous<br>Bounded unreliable<br>Clique | Asynchronous<br>Reliable clique                       | Disconnected                     |
|   | Synchronous<br>Reliable connected<br>Unknown diameter | Can loose all messages           |

# Feasibility: Leader Election



# Models of distributed systems and related faults

## Reading material

► Fred. B. Schneider

“What good are models and what models are good?”

DISTRIBUTED SYSTEMS (2ND ED), SAPE MULLENDER (ED.), ADDISON-WESLEY 1993

