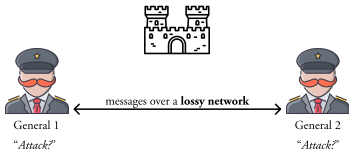


Logical models

Generals problem



General 1	General 2	Outcome
<i>Do not attack</i>	<i>Do not attack</i>	<i>Nothing happens</i>
<i>Attack</i>	<i>Do not attack</i>	<i>General 1 defeated</i>
<i>Do not attack</i>	<i>Attack</i>	<i>General 2 defeated</i>
<i>Attack</i>	<i>Attack</i>	<i>Fortress captured</i>

Goal:

General 1 attacks
if and only if
General 2 attacks

Generals problem



General 1



General 2

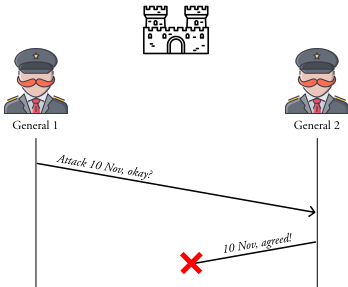
No common knowledge:

only way of knowing something
is to communicate it

Goal:

General 1 attacks
if and only if
General 2 attacks

Generals problem



Goal:

General 1 attacks
if and only if
General 2 attacks

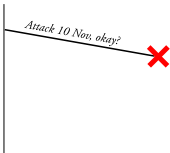
Generals problem



General 1



General 2



Goal:

General 1 attacks
if and only if
General 2 attacks

Option 1:

I attend regularly, but I am not quite
the top management, probably
that General 1 actually knows

Option 2:

I am quite a member of General 1
at least I am not...



General 1



Option 1:

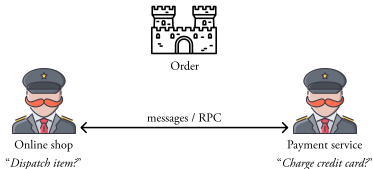
I am prepared to attack... but I am not quite
out of range to reach your facility
and General 1's secondary base.

Option 2:

I am prepared to attack
your facility.

General 2

The two generals problem in practice



Online shop	Payment service	Outcome
<i>Do not dispatch</i>	<i>Do not charge</i>	<i>Nothing happens</i>
<i>Dispatch</i>	<i>Do not charge</i>	<i>Shop loses</i>
<i>Do not dispatch</i>	<i>Charge</i>	<i>Customer unhappy</i>
<i>Dispatch</i>	<i>Charge</i>	<i>Everyone happy</i>

Goal:

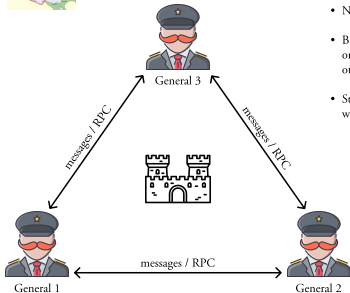
Online shop dispatches
if and only if
Payment service charges

Byzantine generals problem

Byzantine

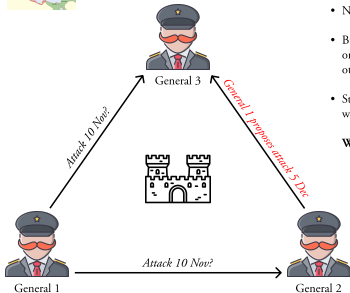


Byzantine generals problem



- Network is now reliable
- But generals are not: one general is **malicious**, other generals don't know who it is
- Still, honest generals would like to agree on a plan

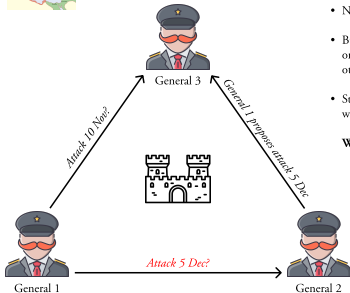
Byzantine generals problem



- Network is now reliable
- But generals are not: one general is **malicious**, other generals don't know who it is
- Still, honest generals would like to agree on a plan

Who can trust whom?

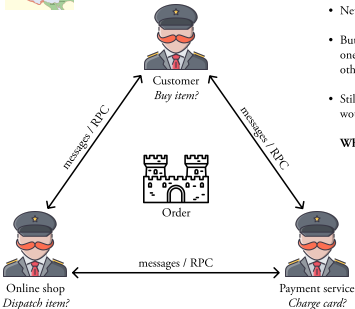
Byzantine generals problem



- Network is now reliable
- But generals are not: one general is **malicious**, other generals don't know who it is
- Still, honest generals would like to agree on a plan

Who can trust whom?

Byzantine generals problem in practice!



- Network is now reliable
- But generals are not: one general is **malicious**, other generals don't know who it is
- Still, honest generals would like to agree on a plan

Who can trust whom?

Faulty behaviours

We have seen two thought experiments:

- **Two generals problem:** a model of *faulty networks*
- **Byzantine generals problem:** a model of *faulty node behaviour*

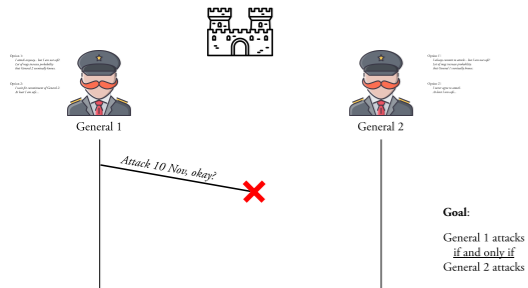
In real systems, both nodes and networks may be faulty!

A generic (failure-oriented) **logical model** of a distributed system may capture:

- *Network behaviour* (e.g. message loss, reordering, duplication, ...)
- *Node behaviour* (e.g. crashes, malicious behaviours, ...)
- *Timing guarantees* (e.g. latency, execution time, clock drifting, ...)

Logical models

Generals problem

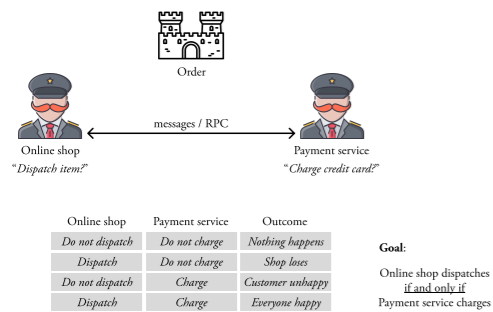


Network behaviour

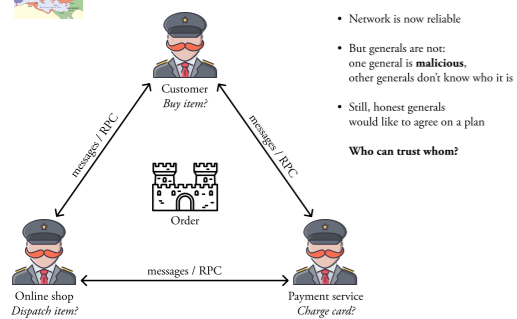
Node behaviour

Timing guarantees

The two generals problem in practice



Byzantine generals problem in practice!



- Network is now reliable
 - But generals are not: one general is **malicious**, other generals don't know who it is
 - Still, honest generals would like to agree on a plan
- Who can trust whom?

Faulty behaviours

We have seen two thought experiments:

- **Two generals problem:** a model of *faulty networks*
- **Byzantine generals problem:** a model of *faulty node behaviour*

In real systems, both nodes and networks may be faulty!

A generic (failure-oriented) **logical model** of a distributed system may capture:

- *Network behaviour* (e.g. message loss, reordering, duplication, ...)
- *Node behaviour* (e.g. crashes, malicious behaviours, ...)
- *Timing guarantees* (e.g. latency, execution time, clock drifting, ...)

Network behaviour

We start by assuming **bidirectional point-to-point** communication between nodes.
So network is represented as an *undirected connected graph*.



COMPLETENESS \rightarrow ∞



Network behaviour



We start by assuming **bidirectional point-to-point** communication between nodes.
So network is represented as an *undirected connected graph*.



Links (channels) between nodes can have at least 4 possible behaviours:

- **Reliable & FIFO:**
A message is received iff it is sent. Messages are received in the same order they are sent.
- **Reliable & non-FIFO:**
A message is received iff it is sent. Messages can be reordered.
- **Fair-loss:**
Messages can be lost, duplicated, reordered. If one keeps retrying, a message eventually gets through.
- **Adversarial:**
A malicious adversary may interfere with messages (eavesdrop, modify, drop, spoof, replay).

- **Fair-loss:**

Messages can be lost

Use of network

resources is minimized

through RTT

- **Adversarial:**

A malicious advers

- **Reliable & non-FIFO**

A message is received in the order it was sent

The message is received in the order it was sent

- **Fair-loss:**

Messages can be lost

The message is received in the order it was sent

- **Reliable & FIFO:**
A message is received in the order it was sent.

- **Reliable & non-FIFO:**
A message is received in an order that is not necessarily the order it was sent.

Node behaviour

Each node executes a *specified algorithm*,
and usually has a *unique identifier* (execution of the algorithm may depend on this).



Node behaviour



Each node executes a *specified algorithm*, and usually has a *unique identifier* (execution of the algorithm may depend on this).

Failure to correctly execute the specified algorithm can result in at least 3 possible behaviours:

- **Crash-stop:**
A node may crash at any moment. After crashing, it stops executing forever.
- **Crash-recovery:**
A node may crash at any moment, losing its in-memory state. It may resume execution some time later, in which case it can recover data stored on disk.
- **Adversarial / Byzantine:**
A node may deviate from the specified algorithm and do anything, including crashing or any malicious action (blocking of messages, replaying, forging, ...).

A node that is not faulty is called **correct**.

Timing guarantees

We may assume one of the following timing guarantees (*):

Synchronous

Asynchronous

(*) in different areas of Computer Science the terms *synchronous* / *asynchronous* have different meanings.

Synchronous

- **every device is present/active**
- **each node requires a time quantum**
- **deficiency: if one device fails, all other will stop too**

Message

- **receiver buffer needed**
- **sender/receiver and processor only involved by request**
- **possible: processing independent from send/receive time**

Asynchronous

Advantages

- easy to use & implement
- needs no special hardware
- independent timing controls

Disadvantages

- less parallelism

Partially synchronous

• Synchronous protocol and not all CPU are active/ready for each iteration

• Good if only one processor used

! Synchronisation usually required to coordinate CPU loads

Violations of timing guarantees in practice

Linux benchmarks measure the time for increasingly larger data sets

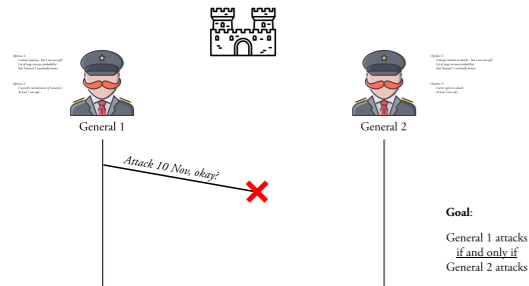
- <http://www.linux-benchmark.com/>
- Computer Associates
- Florida State University

Linux performance regression tests, the process was eventually stopped due to:

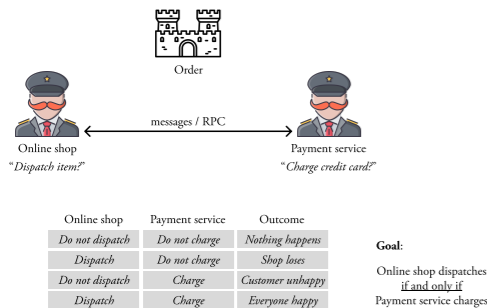
- Unpredictable scheduling causing priority inversion
- High thread-to-pollup utilization
- Page faults, cache thrashing, ...

Logical models

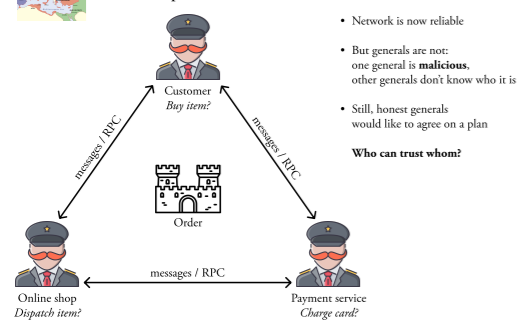
Generals problem



The two generals problem in practice



Byzantine generals problem in practice!



Faulty behaviours

We have seen two thought experiments:

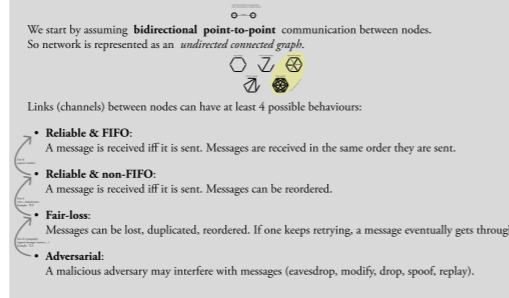
- **Two generals problem:** a model of *faulty networks*
- **Byzantine generals problem:** a model of *faulty node behaviour*

In real systems, both nodes and networks may be faulty!

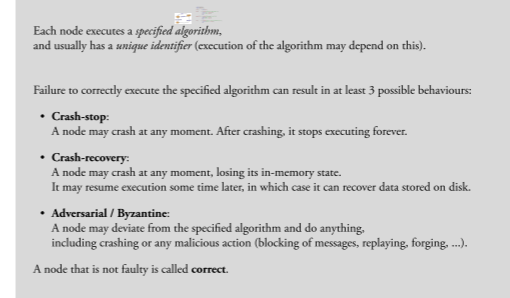
A generic (failure-oriented) **logical model** of a distributed system may capture:

- *Network behaviour* (e.g. message loss, reordering, duplication, ...)
- *Node behaviour* (e.g. crashes, malicious behaviours, ...)
- *Timing guarantees* (e.g. latency, execution time, clock drifting, ...)

Network behaviour



Node behaviour



Timing guarantees



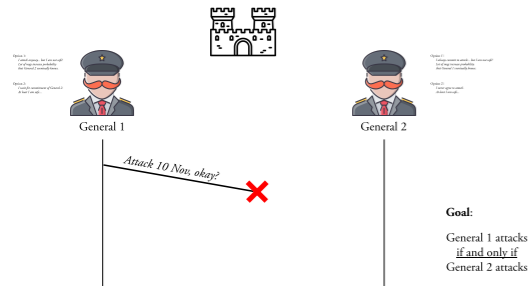
A logical model

e.g. reliable network + crash-stop + synchronous syst.
vs adversary network + Byzantine + asynchronous syst.

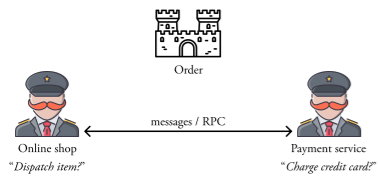
If your assumptions are wrong, all bets are off!

Logical models

Generals problem



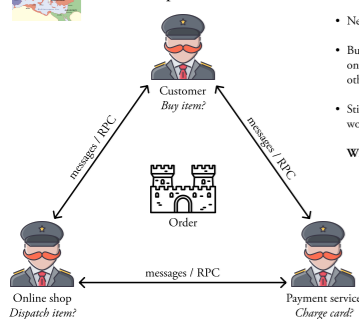
The two generals problem in practice



Online shop	Payment service	Outcome
Do not dispatch	Do not charge	Nothing happens
Dispatch	Do not charge	Shop loses
Do not dispatch	Charge	Customer unhappy
Dispatch	Charge	Everyone happy

Goal: Online shop dispatches if and only if Payment service charges

Byzantine generals problem in practice!



- Network is now reliable
 - But generals are not: one general is **malicious**, other generals don't know who it is
 - Still, honest generals would like to agree on a plan
- Who can trust whom?

Faulty behaviours

We have seen two thought experiments:

- **Two generals problem:** a model of *faulty networks*
- **Byzantine generals problem:** a model of *faulty node behaviour*

In real systems, both nodes and networks may be faulty!

A generic (failure-oriented) **logical model** of a distributed system may capture:

- *Network behaviour* (e.g. message loss, reordering, duplication, ...)
- *Node behaviour* (e.g. crashes, malicious behaviours, ...)
- *Timing guarantees* (e.g. latency, execution time, clock drifting, ...)

Network behaviour

We start by assuming **bidirectional point-to-point** communication between nodes. So network is represented as an *undirected connected graph*.



Links (channels) between nodes can have at least 4 possible behaviours:

- **Reliable & FIFO:** A message is received iff it is sent. Messages are received in the same order they are sent.
- **Reliable & non-FIFO:** A message is received iff it is sent. Messages can be reordered.
- **Fair-loss:** Messages can be lost, duplicated, reordered. If one keeps retrying, a message eventually gets through.
- **Adversarial:** A malicious adversary may interfere with messages (eavesdrop, modify, drop, spoof, replay).

Node behaviour

Each node executes a *specified algorithm*, and usually has a *unique identifier* (execution of the algorithm may depend on this).

Failure to correctly execute the specified algorithm can result in at least 3 possible behaviours:

- **Crash-stop:** A node may crash at any moment. After crashing, it stops executing forever.
 - **Crash-recovery:** A node may crash at any moment, losing its in-memory state. It may resume execution some time later, in which case it can recover data stored on disk.
 - **Adversarial / Byzantine:** A node may deviate from the specified algorithm and do anything, including crashing or any malicious action (blocking of messages, replaying, forging, ...).
- A node that is not faulty is called **correct**.

Timing guarantees

We may assume one of the following timing guarantees (*):

Synchronous **Asynchronous** **Partially synchronous**

Violations of timing guarantees in practice

Violations of timing guarantees in practice

(* In different areas of Computer Science the terms *synchronous / asynchronous* have different meanings.)

A logical model

e.g. reliable network + crash-stop + synchronous syst.
vs adversary network + Byzantine + asynchronous syst.

If your assumptions are wrong, all bets are off!

A good practice: try to avoid single points of failure.

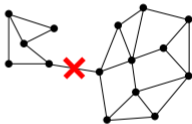
Network partitioning

Network partitioning

Network / node faults may sometimes lead to a **partitioning of the network** (a.k.a. “split brain”), where some nodes may become disconnected from other nodes, and may assume the latter have crashed.

Even if this fault were temporary, it could lead to:

- *inconsistency* (e.g. of replicated data)
- *unavailability* (e.g. of a resource or a service)



CAP Theorem

A system can be either (strongly) **consistent (C)** or **available (A)** in presence of **network partition (P)**.

consistent (C)

Consistent

Consistent

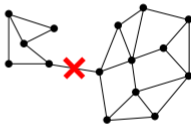
available (A)

Network partitioning

Network / node faults may sometimes lead to a **partitioning of the network** (a.k.a. “split brain”), where some nodes may become disconnected from other nodes, and may assume the latter have crashed.

Even if this fault were temporary, it could lead to:

- *inconsistency* (e.g. of replicated data)
- *unavailability* (e.g. of a resource or a service)



CAP Theorem

A system can be either (strongly) **consistent (C)** or **available (A)** in presence of **network partition (P)**.

C_{onsistent} + **P**_{artition-tolerant}

A_{vailable} + **P**_{artition-tolerant}

Causation + Prediction

Relationships between past & equilibrium climate states.

Example

- 1870
- 1990-2010

A_{available} + P_{ermissions-demand}

Successful response to information technology.

Example

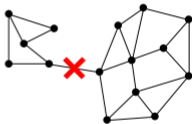
- 1987
- 1990-1992

Network partitioning

Network / node faults may sometimes lead to a **partitioning of the network** (a.k.a. “split brain”), where some nodes may become disconnected from other nodes, and may assume the latter have crashed.

Even if this fault were temporary, it could lead to:

- *inconsistency* (e.g. of replicated data)
- *unavailability* (e.g. of a resource or a service)



CAP Theorem

A system can be either (strongly) **consistent (C)** or **available (A)** in presence of **network partition (P)**.

Consistent + **P**artition-tolerant

Available + **P**artition-tolerant

Eventually **C**onsistent + **A**vailable + **P**artition-tolerant

Eventually **C**onsistent + **A**vailable + **P**artition-tolerant

If an update is made to a given resource,
that eventually (but no-guarantee how long this might take),
from all nodes, all access to that resource will return the last updated value of that resource.

Since due to possible concurrent updates, this requires a mechanism for conflict resolution.

Examples:

- Blockchain
- Amazon Web Services

Fault tolerance

Fault: *some part* of the system is not working correctly, e.g.

- node fault (e.g. crash-stop, crash-recovery, Byzantine deviation, ...)
- network fault (e.g. dropping or significantly delaying msgs, evedropping or forgery of msgs, ...)

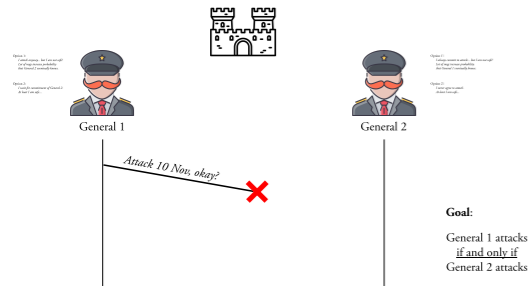
Failure: *system as a whole* is not functioning correctly (e.g. it is temporarily unavailable)

Fault tolerance: system as a whole continues working, despite faults (up to some maximum number of faults).

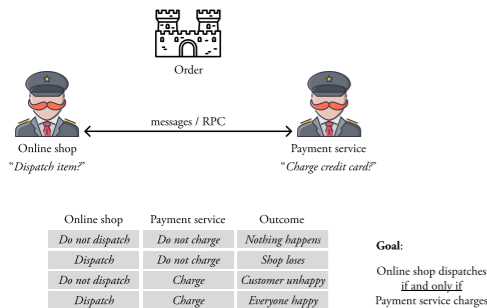
Usually requires mechanisms for *fault detection*.



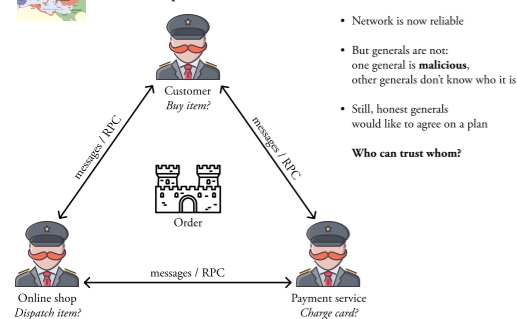
Generals problem



The two generals problem in practice



Byzantine generals problem in practice!



Faulty behaviours

We have seen two thought experiments:

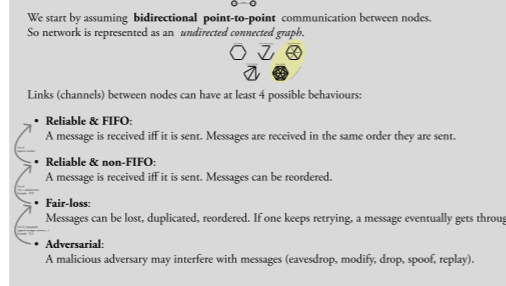
- Two generals problem:** a model of *faulty networks*
- Byzantine generals problem:** a model of *faulty node behaviour*

In real systems, both nodes and networks may be faulty!

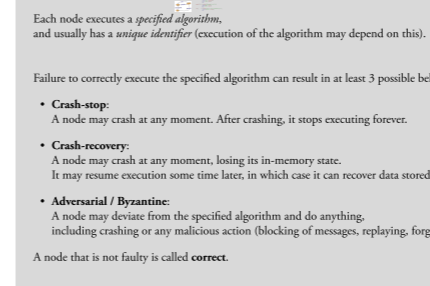
A generic (failure-oriented) **logical model** of a distributed system may capture:

- Network behaviour** (e.g. message loss, reordering, duplication, ...)
- Node behaviour** (e.g. crashes, malicious behaviours, ...)
- Timing guarantees** (e.g. latency, execution time, clock drifting, ...)

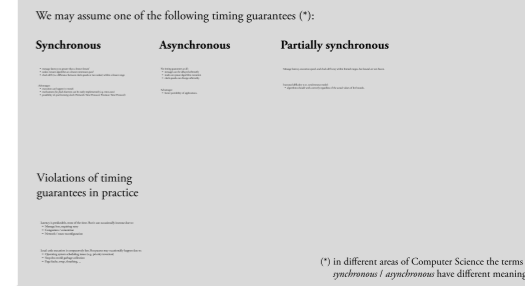
Network behaviour



Node behaviour



Timing guarantees



A logical model

e.g. reliable network + crash-stop + synchronous syst.
vs adversary network + Byzantine + asynchronous syst.

If your assumptions are wrong, all bets are off!

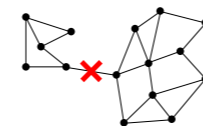
A good practice: try to avoid single points of failure.

Network partitioning

Network / node faults may sometimes lead to a **partitioning of the network** (a.k.a. "split brain"), where some nodes may become disconnected from other nodes, and may assume the latter have crashed.

Even if this fault were temporary, it could lead to:

- inconsistency** (e.g. of replicated data)
- unavailability** (e.g. of a resource or a service)



CAP Theorem

A system can be either (strongly) **consistent (C)** or **available (A)** in presence of **network partition (P)**.



Fault tolerance

Fault: *some part* of the system is not working correctly, e.g.

- node fault (e.g. crash-stop, crash-recovery, Byzantine deviation, ...)
- network fault (e.g. dropping or significantly delaying msgs, evedropping or forgery of msgs, ...)

Failure: *system as a whole* is not functioning correctly (e.g. it is temporarily unavailable)

Fault tolerance: system as a whole continues working, despite faults (up to some maximum number of faults).

Usually requires mechanisms for **fault detection**.

