**Distributed Systems** 

## Byzantine Fault Tolerance and Consensus

Dragos-Adrian Seredinschi Ipd.epfl.ch

## Problem



#### General goal: Run a distributed algorithm

# Problem



#### General goal: Run a distributed algorithm

## Recasting the problem



# Recasting the problem

Application requirements:

- High-Availability (give a reply to a request)
- Reliability (give correct replies)

Boils down to fault-tolerance



#### Solution Fault-tolerance basic techniques: Client • Agreement = Consensus REPLY Application REQUEST **State Machine Replication** Distributed algorithm

In the following we will see...

PBFT  $\bullet$ 

Replication =

Seminal algorithm for • **Byzantine Fault Tolerance** 

#### PBFT Practical Byzantine Fault Tolerance OSDI'99

![](_page_6_Picture_1.jpeg)

Miguel Castro

![](_page_6_Picture_3.jpeg)

Barbara Liskov

# Modules

![](_page_7_Figure_1.jpeg)

# Modules

![](_page_8_Figure_1.jpeg)

# Overview

#### PBFT:

#### · System model

- slightly different from what we've seen so far

- SMR
- Consensus

### System model Processes

Three types of processes in this algorithm:

• Clients

![](_page_10_Figure_3.jpeg)

- n replicas
  - one of them is primary
  - others are <u>backup</u>

![](_page_10_Picture_7.jpeg)

## System model Failure model

- Arbitrary (Byzantine) faults
- Clients:
  - Any client can be faulty
- Replicas:
  - n = 3f + 1
  - f faulty (upper bound)

e

9

n=4 f=1

n=7

f=2

e

## System model Network & crypto

- Assume perfect links
- Direct links between any two processes
- For messages:
  - Public-key signatures, message authentication codes
  - Avoid spoofing, replays, corruption
- Clients are authenticated
  - Can revoke access to faulty clients

# Overview

#### PBFT:

• System model

- slightly different from what we've seen so far

#### · SMR

• Consensus

# State Machine Replication

- A fault-tolerance technique
- Basic ideas:
  - Application = state machine

![](_page_14_Figure_4.jpeg)

- Run the application on multiple processes
- Each processes is a faithful replica of the application
- Note: We can ignore the primary/backup distinction in this example

## SMR in a nutshell

![](_page_15_Figure_1.jpeg)

![](_page_16_Figure_0.jpeg)

![](_page_17_Figure_0.jpeg)

![](_page_18_Figure_0.jpeg)

![](_page_19_Figure_0.jpeg)

![](_page_20_Figure_0.jpeg)

![](_page_21_Figure_0.jpeg)

## SMR Requirements

- Avoid diverging states
- All replicas must:
  - 1. Start in the same state
  - 2. Execute the same sequence of operations
  - Use only provided operation (+parameters), thus avoid non-determinism

![](_page_23_Figure_0.jpeg)

- Application = a distributed file system
   Network File System (NFS)
- Operations = write to a file, delete, etc.
- Primary/backup distinction is relevant

NFS
State Machine Replication (SMR)
Consensus
Perfect Links
Channels

![](_page_24_Picture_5.jpeg)

- Application = a distributed file system
   Network File System (NFS)
- Operations = write to a file, delete, etc.
- Primary/backup distinction is relevant

	NFS
	State Machine Replication (SMR)
	Consensus
-	Perfect Links
	Channels

![](_page_25_Figure_5.jpeg)

![](_page_26_Figure_1.jpeg)

![](_page_27_Figure_1.jpeg)

![](_page_28_Figure_0.jpeg)

![](_page_29_Figure_0.jpeg)

![](_page_30_Figure_0.jpeg)

# Overview

#### PBFT:

• System model

- slightly different from what we've seen so far

- SMR
- · Consensus

# Consensus

- The core for many algorithms, including:
- TRB, Group membership, View synchronous b-cast, State machine replication

#### Traditionally

 Processes propose values

Instance

Agree on a proposed value

#### In PBFT:

- Clients propose <u>request</u>
- Primary multicasts one request to backup replicas
- Replicas accept the request

nstance

- We'll assume one client
  - <u>Proposals = requests</u> for application operations
- Assume:
  - n = 4, f = 1
  - The faulty replica does not cooperate
- Concurrent requests:
  - <u>Consensus</u> to agree on a sequential execution of requests

![](_page_33_Picture_8.jpeg)

Algorithm ideas:

- Client sends requests to the primary replica
- Execute a sequence of consensus instances:
  - Each instance is dedicated to a request
  - Instances (and therefore requests) are sequentially ordered by the primary
  - Backup replicas adopt requests from the primary in the imposed order

Properties: Validity, Agreement, Termination, Integrity

![](_page_35_Figure_1.jpeg)

![](_page_36_Figure_1.jpeg)

#### A three-phase protocol

![](_page_37_Figure_1.jpeg)

![](_page_38_Figure_1.jpeg)

![](_page_39_Figure_1.jpeg)

![](_page_40_Figure_1.jpeg)

![](_page_41_Figure_1.jpeg)

![](_page_42_Picture_0.jpeg)

What if the primary is faulty, e.g. does not multicast the request to the backups?

- View change protocol: primary replaced by one of the backups
- Idea:
  - Replicas are numbered 1 ... n
  - In view v, the replica p is the primary, where
    p = v mod n

# Practical BFT

"Reasonable overhead"

- Does not assume synchrony
- Some clever optimizations:
  - MD5 replaces digital signatures
  - Message digests
  - Read-only requests, tentative execution

![](_page_43_Picture_7.jpeg)

# Further reading

- Castro, M., & Liskov, B. (1999). Practical Byzantine fault tolerance. OSDI, (February), 1–14. Available at: <u>http://dl.acm.org/citation.cfm?id=296806.296824</u>
- Castro, M. (2011). Practical Consensus. Microsoft Research Cambridge. Available at: <a href="http://msrvideo.vo.msecnd.net/rmcvideos/167097/dl/167097.pdf">http://msrvideo.vo.msecnd.net/rmcvideos/167097/dl/167097.pdf</a>