



Arhitecturi Paralele

Consistency and Consensus

Prof. Florin Pop
As. Drd. Ing. Cristian Chilipirea
cristian.chilipirea@cs.pub.ro

Elemente preluate din cursul Prof. Ciprian Dobre



FACULTATEA DE
**AUTOMATICĂ ȘI
CALCULATOARE**





Cum se măsoară viteza unui proces și ce ne spune?



În cât timp va executa următorul program?

```
int main() {  
    → int i;  
  
    → for (i=0; i < 1000 * 1000 * 1000; i++);  
  
    → return 0;  
}
```

În cât timp va executa următorul program?

```
int main() {  
    → int i;  
    → for (i=0; i < 1000 * 1000 * 1000; i++);  
    → return 0;  
}
```

```
real    0m1.783s  
user    0m1.766s  
sys     0m0.000s
```

Time - precizie

Scientific Notation

Term	Symbol	Multiple	Sci Not ⁿ
Terra	T	1 000 000 000 000	$\times 10^{12}$
Giga	G	1 000 000 000	$\times 10^9$
Mega	M	1 000 000	$\times 10^6$
kilo	k	1 000	$\times 10^3$
Units	Eg. meters	1	$\times 10^0$
milli	m	1 / 1 000	$\times 10^{-3}$
micro	μ	1 / 1 000 000	$\times 10^{-6}$
nano	n	1 / 1 000 000 000	$\times 10^{-9}$
pico	p	1 / 1000 000 000 000	$\times 10^{-12}$





Timp transmite pachet?



Timp transmiere pachet?

```
root@Nyx:/mnt/d/Dropbox/backupServer/apd-homework/rezults# ping 8.8.8.8
PING 8.8.8.8 (8.8.8.8) 56(84) bytes of data.
64 bytes from 8.8.8.8: icmp_seq=1 ttl=120 time=15.4 ms
64 bytes from 8.8.8.8: icmp_seq=2 ttl=120 time=13.0 ms
64 bytes from 8.8.8.8: icmp_seq=3 ttl=120 time=15.3 ms
64 bytes from 8.8.8.8: icmp_seq=4 ttl=120 time=14.5 ms
64 bytes from 8.8.8.8: icmp_seq=5 ttl=120 time=32.5 ms
64 bytes from 8.8.8.8: icmp_seq=6 ttl=120 time=15.2 ms
```


Sincronizarea timpilor între 2 sisteme

Probabilistic clock synchronization

Flaviu Cristian

IBM Almaden Research Center, 650 Harry Road, San Jose, CA 95120, USA



Flaviu Cristian is a computer scientist at the IBM Almaden Research Center in San Jose, California. He received his PhD from the University of Grenoble, France, in 1979. After carrying out research in operating systems and programming methodology in France, and working on the specification, design, and verification of fault-tolerant programs in England, he joined IBM in 1982. Since then he has worked in the area of fault-tolerant distributed protocols and

systems. He has participated in the design and implementation of a highly available system prototype at the Almaden Research Center and has reviewed and consulted for several fault-tolerant distributed system designs, both in Europe and in the American divisions of IBM. He is now a technical leader in the design of a new U.S. Air Traffic Control System which must satisfy very stringent availability requirements.

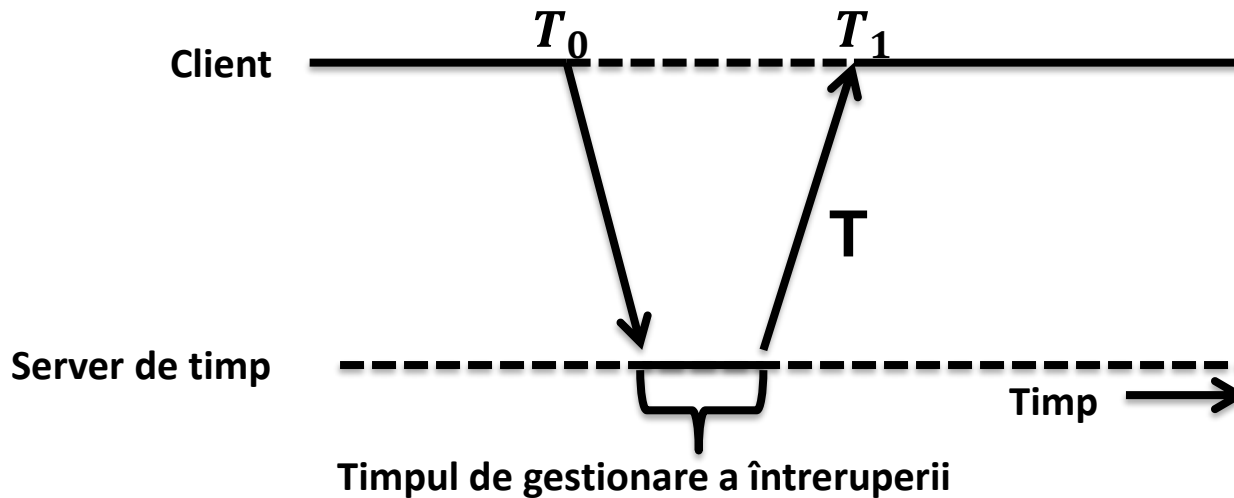
Abstract. A probabilistic method is proposed for reading remote clocks in distributed systems subject to unbounded random communication delays.

in some given maximum derivation from a time reference external to the system. *Internal* clock synchronization keeps processor clocks within some maximum relative deviation of each other. Externally synchronized clocks are also internally synchronized. The converse is not true: as time passes internally synchronized clocks can drift arbitrarily far from external time.

Clock synchronization is needed in many distributed systems. Internal clock synchronization enables one to measure the duration of distributed activities that start on one processor and terminate on another processor and to totally order distributed events in a manner that closely approximates their real time precedence. To allow exchange of information about the timing of events with other systems and users, many systems require external clock synchronization. For example external time can be used to record the occurrence of events for later analysis by humans, to instruct a system to take certain actions when certain specified (external) time deadlines occur, and to order the occurrence of related events observed by distinct sys-

Sincronizarea timpilor între 2 sisteme

$$T + \text{RTT}/2$$





Cu ce protocol se sincronizează ceasul calculatoarelor?

Cu ce protocol se sincronizează ceasul calculatoarelor?

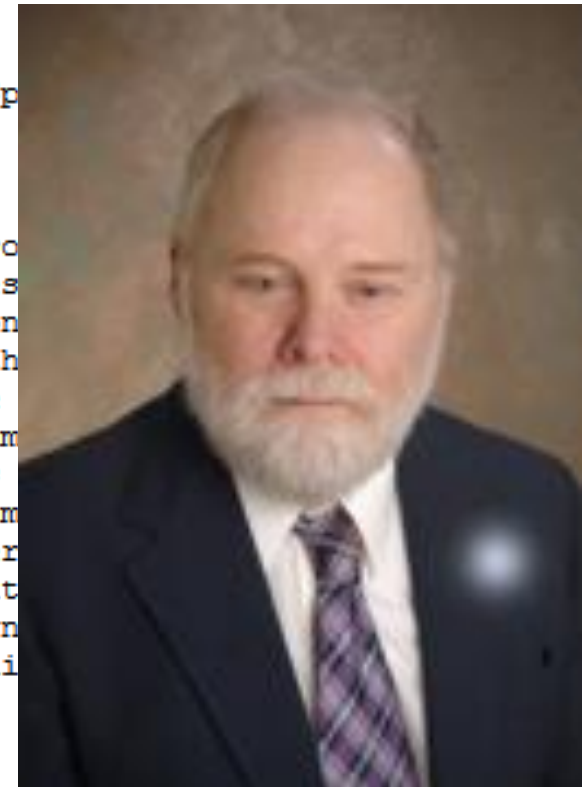
Internet Engineering Task Force (IETF)
Request for Comments: 5905
Obsoletes: 1305, 4330
Category: Standards Track
ISSN: 2070-1721

D. Mills
U. Delaware
J. Martin, Ed.
ISC
J. Burbank
W. Kasch
JHU/APL
June 2010

Network Time Protocol Version 4: Protocol and Algorithms Specification

Abstract

The Network Time Protocol (NTP) is widely used to synchronize computer clocks in the Internet. This document describes NTP version 4 (NTPv4), which is backwards compatible with NTP version 3 described in RFC 1305, as well as previous versions of the protocol. NTPv4 includes a modified protocol header to accommodate Protocol version 6 address family. NTPv4 includes fundamental improvements in the mitigation and discipline algorithms, achieving the potential accuracy to the tens of microseconds with most workstations and fast LANs. It includes a dynamic server selection scheme, so that in many cases, specific server configuration is not required. It corrects certain errors in the NTPv3 design and implementation and includes an optional extension mechanism.





Intr-un sistem distribuit cum stabilim ordinea operațiilor?



Modelarea unui sistem distribuit?



Modelarea unui sistem simplu/secvențial?

Modelarea unui sistem simplu/secvențial?

First Draft of a Report
on the EDVAC

by

John von Neumann





Modelarea unui sistem simplu/secvențial?

Un set de stări.

Operațiile sau acțiunile modifică starea



Modelarea unui sistem distribuit?

Modelarea unui sistem distribuit?

Programming
Techniques

S. L. Graham, R. L. Rivest
Editors

Communicating Sequential Processes

C.A.R. Hoare
The Queen's University
Belfast, Northern Ireland

This paper suggests that input and output are basic primitives of programming and that parallel composition of communicating sequential processes is a fundamental program structuring method. When combined with a development of Dijkstra's guarded command, these concepts are surprisingly versatile. Their use is illustrated by sample solutions of a variety of familiar programming exercises.

Key Words and Phrases: programming, programming languages, programming primitives, program structures, parallel programming, concurrency, input, output, guarded commands, nondeterminacy, coroutines, procedures, multiple entries, multiple exits, classes, data representations, recursion, conditional critical regions, monitors, iterative arrays

CR Categories: 4.20, 4.22, 4.32

grams, three basic constructs have received widespread recognition and use: A repetitive construct (e.g. the **while** loop), an alternative construct (e.g. the conditional **if..then..else**), and normal sequential program composition (often denoted by a semicolon). Less agreement has been reached about the design of other important program structures, and many suggestions have been made: Subroutines (Fortran), procedures (Algol 60 [15]), entries (PL/I), coroutines (UNIX [17]), classes (SIMULA 67 [5]), processes and monitors (CLU [13]), forms (ALPHA [13]).

The traditional stored-program computer has been designed primarily for single sequential programs. The need for speed has led to the introduction of parallelism. An attempt has been made to design a computer for a programmer, either by having function units of the CDC (in an I/O control package generating system). However, technology suggests that a computer should be structured from a number of processors (each with its own powerful, capacious, reliable machine which is disguised).

In order to use such a computer for a task, the component processes must communicate and to synchronize methods of achieving this. The adopted method of communication is the updating of a common state and many machine code severe problems in the com-





Intr-un sistem distribuit cum stabilim ordinea operațiilor?

Lamport timestamps

Operating Systems

R. Stockton Gaines
Editor

Time, Clocks, and the Ordering of Events in a Distributed System

Leslie Lamport
Massachusetts Computer Associates, Inc.

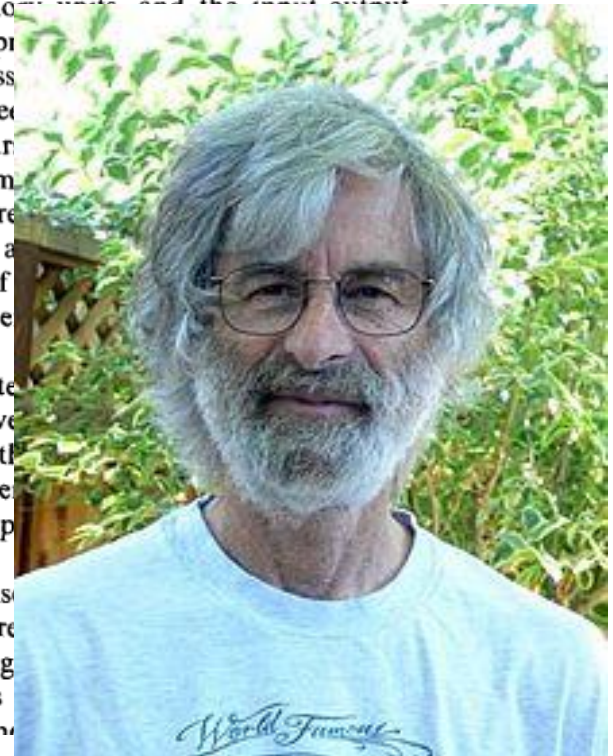
The concept of one event happening before another in a distributed system is examined, and is shown to define a partial ordering of the events. A distributed algorithm is given for synchronizing a system of logical clocks which can be used to totally order the events. The use of the total ordering is illustrated with a method for solving synchronization problems. The algorithm is then specialized for synchronizing physical clocks, and a bound is derived on how far out of synchrony the clocks can become.

A distributed system consists of a collection of distinct processes which are spatially separated, and which communicate with one another by exchanging messages. A network of interconnected computers, such as the ARPANet, is a distributed system. A single computer can also be viewed as a distributed system in which the central control unit, the memory unit, and the input/output channels are separate processes. This is true if the message transmission time is not negligible compared to the time between successive messages.

We will concern our spatially separated components. Our remarks will apply more to a multiprocessing system on a multiprocessor than to problems similar to those of the unpredictable order of events that may occur.

In a distributed system, we say that one of two events “happened before” is the case if the first event is one of the events in the system that often arise because people and its implications.

In this paper, we discuss the algorithm for extending the "happened before" relation of all the events. This mechanism for implementing



Ordine parțială

- **Timp relativ versus timp real**
- **Ceasuri logice**
- **petrecut înainte (\rightarrow)**
 - ◆ dacă **a** și **b** sunt evenimente din același proces și **a** îl precede în timp pe **b**,
 - atunci $a \rightarrow b$
 - ◆ când **a** reprezintă transmiterea unui mesaj de către un proces, iar **b** recepția
— acelui mesaj de către un altul, atunci $a \rightarrow b$
 - ◆ dacă $a \rightarrow b$ și $b \rightarrow c$, atunci $a \rightarrow c$

Ordine parțială

■ Reguli:

- la producerea unui eveniment intern:
 - cl este incrementat
 - valoarea ceasului logic este asociată evenimentului ca **amprentă de timp**
- la transmiterea unui mesaj :
 - incrementează cl al procesului transmitator cu 1
 - actualizează $tt := cl$
- la primirea unui mesaj cu amprenta de timp tt :
 - actualizează $cl := \max(cl, tt) + 1$

Ordine parțială

■ Reguli:

- la producerea unui eveniment intern:
 - **cl** este incrementat
 - valoarea ceasului logic este asociată evenimentului ca **amprentă de timp**
- la transmiterea unui mesaj :
 - incrementează **cl** al procesului transmitator cu 1
 - actualizează **tt** := **cl**
- la primirea unui mesaj cu amprenta de timp **tt**:
 - actualizează **cl** := $\max(\text{cl}, \text{tt}) + 1$

■ Ordonare parțială: $a \rightarrow b \Rightarrow \text{cl}(a) < \text{cl}(b)$

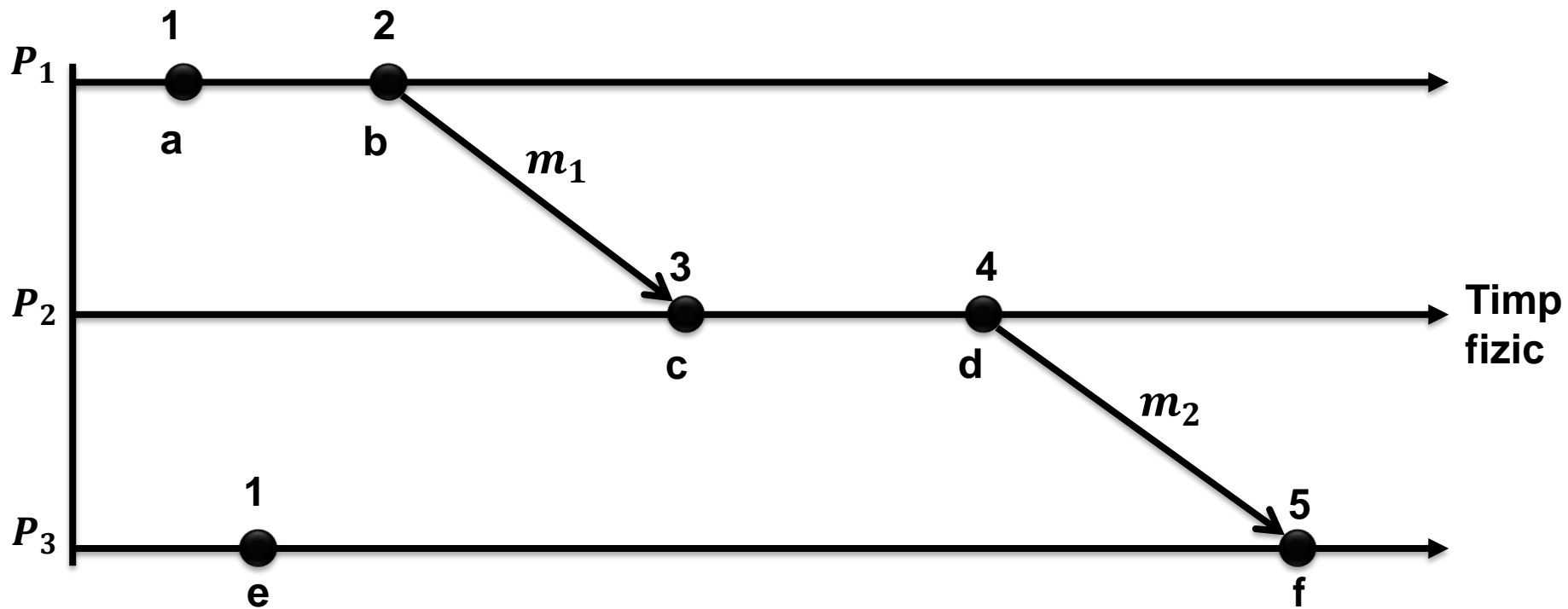
Ordine parțială

■ Reguli:

- la producerea unui eveniment intern:
 - **cl** este incrementat
 - valoarea ceasului logic este asociată evenimentului ca **amprentă de timp**
- la transmiterea unui mesaj :
 - incrementează **cl** al procesului transmitator cu 1
 - actualizează **tt** := **cl**
- la primirea unui mesaj cu amprenta de timp **tt**:
 - actualizează **cl** := $\max(\text{cl}, \text{tt}) + 1$

■ Ordonare parțială: $a \rightarrow b \Rightarrow \text{cl}(a) < \text{cl}(b)$

Exemplu



Ceasuri logice vectoriale

Cu soluția Lamport:

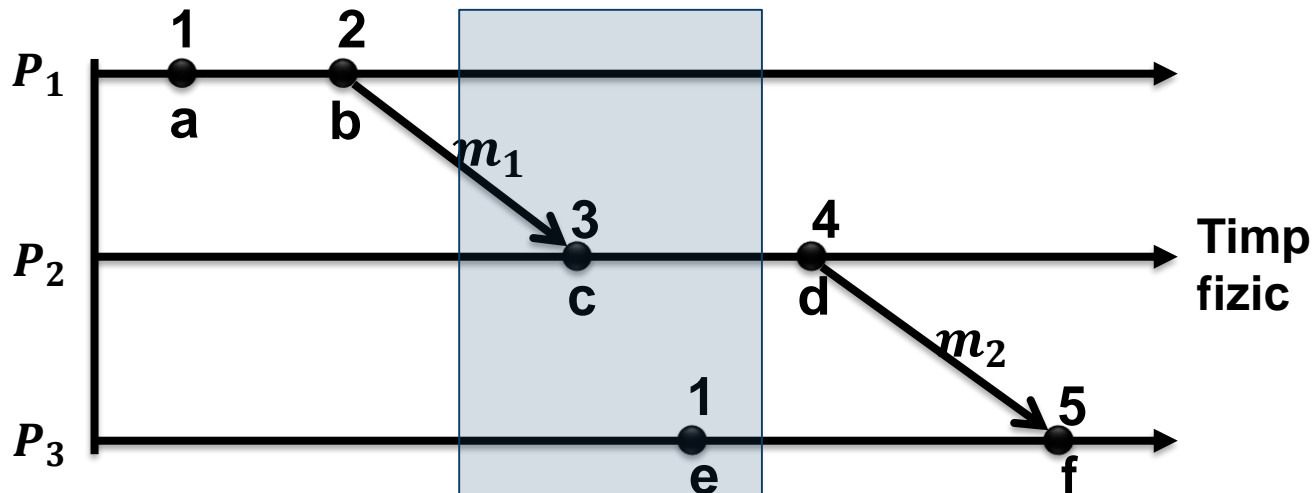
$e \text{ precede } f \Rightarrow \text{amprenta_logică}(e) < \text{amprenta_logică}(f)$

Dar

$\text{amprenta_logică}(e) < \text{amprenta_logică}(f) \not\Rightarrow e \text{ precede } f$

Ex: se poate spune ca **e** precede **c** ?

Soluția: ceasuri logice vectoriale.





Vector Clocks

Timestamps in Message-Passing Systems That Preserve the Partial Ordering

Colin J. Fidge

Department of Computer Science, Australian National University, Canberra, ACT.

ABSTRACT

Timestamping is a common method of totally ordering events in concurrent programs. However, for applications requiring access to the global state, a total ordering is inappropriate. This paper presents algorithms for timestamping events in both synchronous and asynchronous message-passing programs that allow for access to the global state without being herent in a parallel system. The algorithms do not change the consistency of the system and require a central timestamp issuing authority.

Keywords and phrases: concurrent programming, message-passing, timestamps
CR categories: D.1.3

INTRODUCTION

A fundamental problem in concurrent programming is determining the order of events that occurred. An obvious solution is to attach a number representing the order of events to each event. This assumes that each process has a permanent record of the execution of each event. This assumes that each process has a clock, but practical parallel systems, by their very nature, make it difficult to



Vector Clocks

Virtual Time and Global States of Distributed Systems *

Friedemann Mattern [†]

Department of Computer Science, University of Kaiserslautern
D 6750 Kaiserslautern, Germany

Abstract

A distributed system can be characterized by the fact that the global state is distributed and that a common time base does not exist. However, the notion of time is an important concept in every day life of our decentralized “real world” and helps to solve problems like getting a consistent population census or determining the potential causality between events. We argue that a linearly ordered structure of time is not (always) adequate for distributed systems and propose a generalized non-standard model of time which consists of vectors

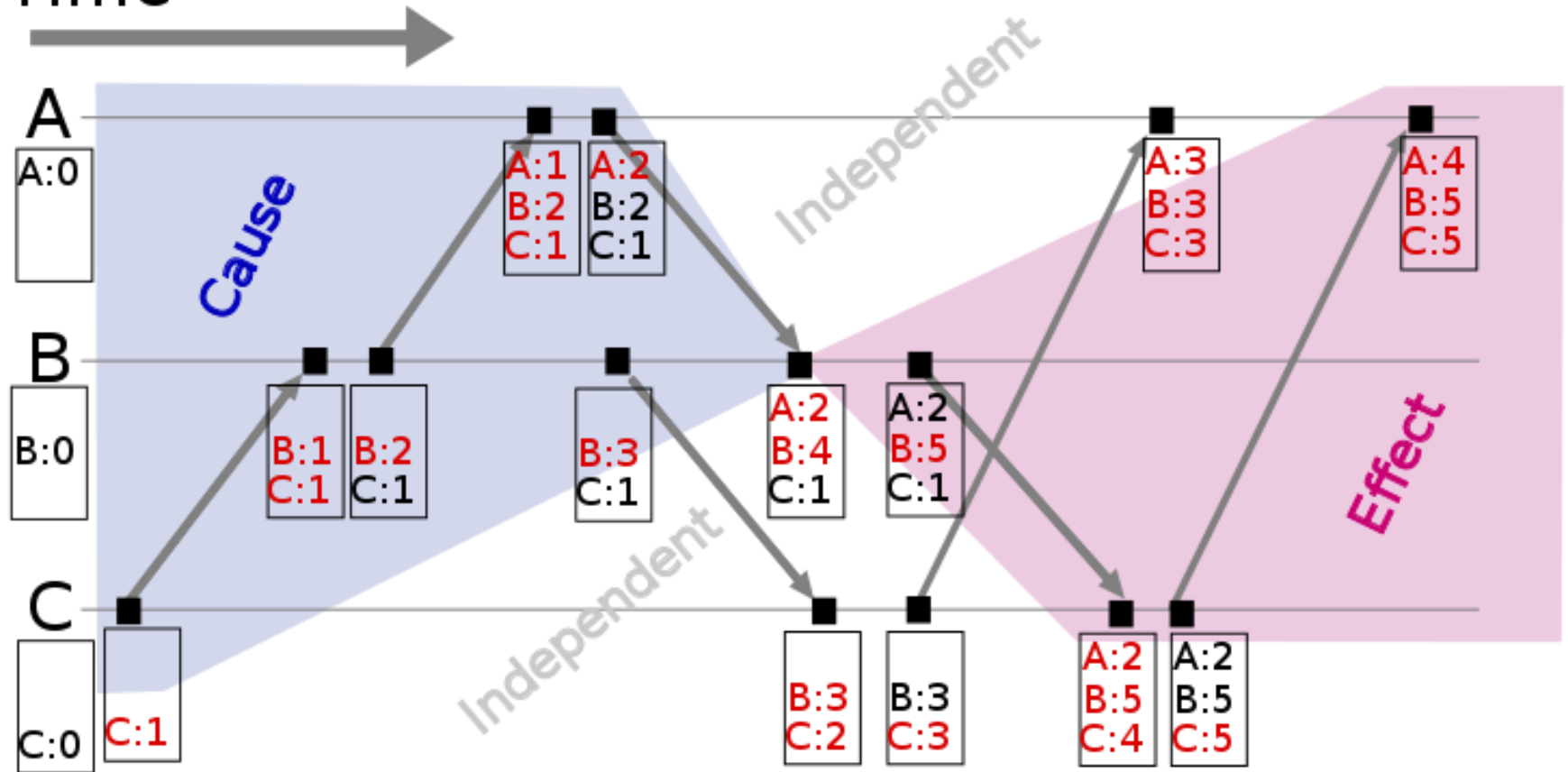
view of an idealized access to all process

The fact that a view of the global state does not exist is the cause of many problems in distributed systems. Consistency, deadlock detection, and concurrency control are difficult to solve in a distributed environment. In a physical centralized environment, the solution of distributed control problems is found to be wrong. In a centralized system,



Vector Clocks

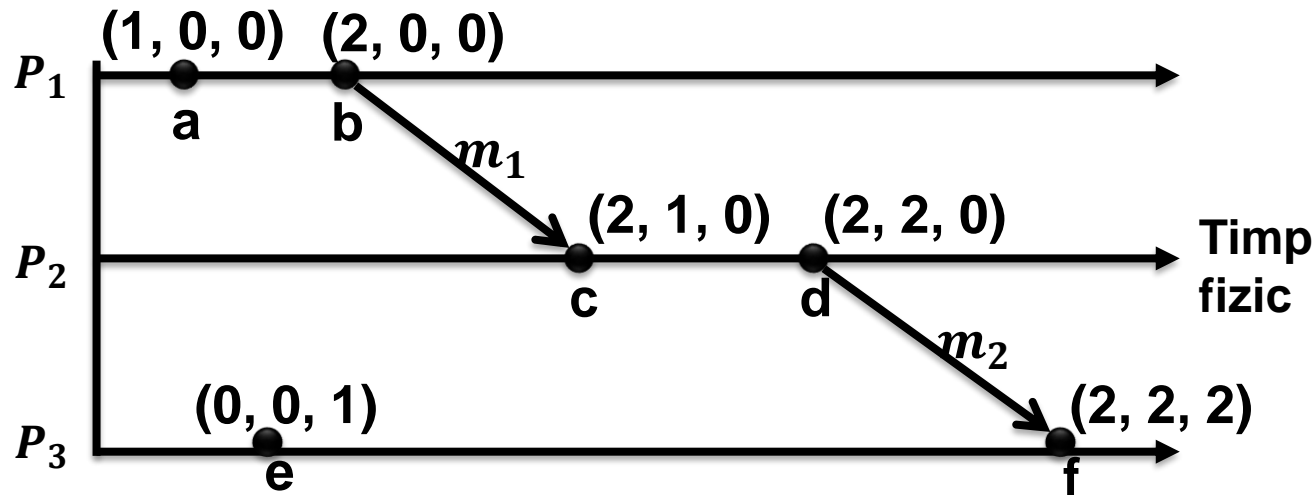
Time



Ceasuri logice vectoriale

- Fiecare P_i are asociat un tablou $V_i[1..n]$ în care:
 - $V_i[i]$ este numărul de **evenimente** produse în procesul P_i
 - $V_i[j]$ este numărul de **evenimente** despre care P_i știe (*a aflat*) că au avut loc la P_j .
- Procesul P_i actualizează V_i la fiecare **eveniment** din P_i
 - e.g. pentru procesorul 3, $(1,2,1,3) \rightarrow (1,2,2,3)$
- Când P_i **transmite** mesajul m (*eveniment send m*):
 - P_i incrementează $V_i[i]$
 - P_i adaugă V_i la m ca vector de amprente de timp curent vt_m
- Când P_j **primește** m și vt_m (*eveniment receive m*):
 - P_j ajustează: $V_j[k] = \max\{V_j[k], vt_m[k]\}$ pentru fiecare k
e.g., P_2 primește un mesaj cu timpul $(3,2,4)$ iar timpul curent al lui P_2 este $(3,4,3)$, atunci P_2 ajustează timpul la $(3,4,4)$
 - P_j incrementează $V_j[j]$ cu 1

Ceasuri logice vectoriale



Aplicarea regulilor ceasurilor logice vectoriale

Reguli:

- $VT_1 = VT_2 \Leftrightarrow VT_1[i] = VT_2[i]$, pentru $i = 1, \dots, N$
- $VT_1 \leq VT_2 \Leftrightarrow VT_1[i] \leq VT_2[i]$, pentru $i = 1, \dots, N$
- $VT_1 < VT_2 \Leftrightarrow VT_1[i] \leq VT_2[i]$, și $VT_1 \neq VT_2$ (de exemplu $(1,2,2) < (1,3,2)$)

Fie $vt(a)$ și $vt(b)$ vectorii de amprente de timp asociați ev. a și b . Atunci:

$vt(a) < vt(b) \Rightarrow$ **evenimentul a precede cauzal b**

$vt(a) \not< vt(b)$ and $vt(a) \not> vt(b)$ and $vt(a) \neq vt(b) \Leftrightarrow$ **evenimentele a și b sunt concurente**

Aplicație: Ordonare Causală Multicast

- Procesele unei colecții P comunică între ele doar prin **mesaje cu difuzare**
- Se cere ca mesajele să respecte **dependența causală**

$$m \rightarrow m' \Rightarrow \text{livrarep}(m) \rightarrow \text{livrarep}(m')$$

Protocolul (*vectori de timp*):

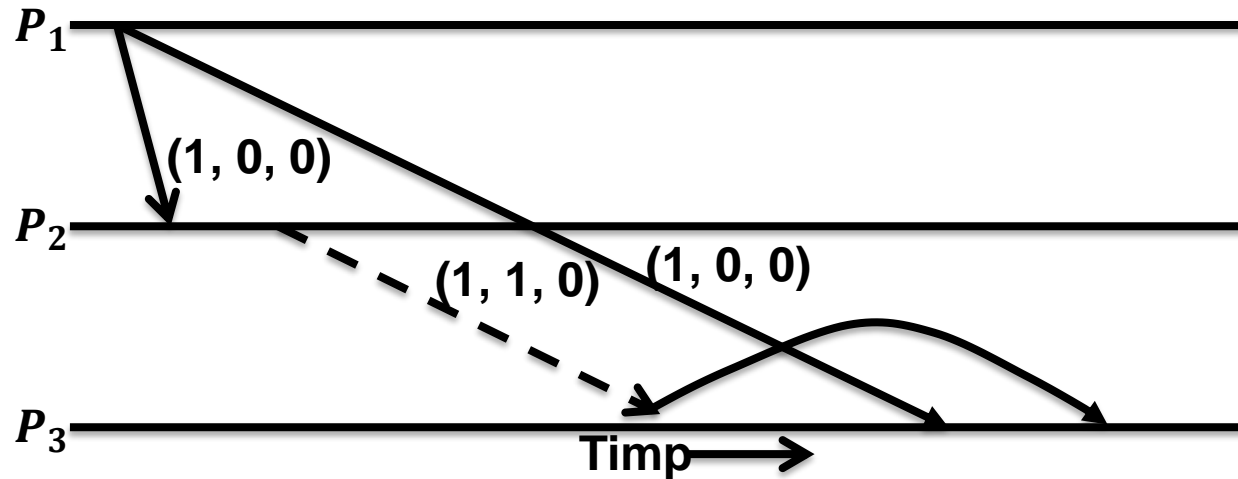
- fiecare proces $P_i = (i = 1..n)$ are asociat un vector $V_i[1..n]$, cu toate elementele inițial 0
- se numara doar operatiile de **transmitere** de mesaje
- $V_i[i]$, este nr ev. **transmitere** de mesaje produse de P_i ;
- $V_i[j]$ este nr ev. transmitere de mesaje **despre care P_i știe** (a aflat) că au avut loc la P_j .
- Procesul P_i actualizează V_i la fiecare eveniment de trimitere sau recepție de mesaj din P_i .



Aplicație: Ordonare Cauzală Multicast

- Când P_s transmite mesajul m :
 - P_s incrementează $V_s[s]$
 - P_s adaugă V_s la m ca vector de amprente de timp curent vt_m
 - **Obs: Pentru ordonare cauzală, incrementarea lui $V_s[s]$ se face doar la transmiterea de mesaje de către s**
 - vt_m spune receptorului câte evenimente (din alte procese) au precedat m și ar putea influența cauzal pe m .
- Când P_d primește mesajul m împreună cu vt_m , mesajul este păstrat într-o coadă de întârziere și este livrat doar dacă:
 - $vt_s[s] = V_d[s] + 1$ (acesta este următorul timestamp pe care d îl așteaptă de la s)
 - $vt_m[k] \leq V_d[k]$ pentru $k \neq s$ (d a văzut toate mesajele ce au fost văzute de s la momentul când a trimis mesajul m)
- Când mesajul este livrat, V_d este actualizat conform regulilor vectorilor de timp:
$$V_d[k] = \max\{V_d[k], vt_m[k]\} \text{ pentru fiecare } k = 1, n.$$

Aplicație: Ordonare Cauzală Multicast



Alte acțiuni ale protocolului la livrarea mesajelor:

- dacă $vt_m[k] > V_d[k]$ pentru un oarecare k atunci se întârzie m
 P_s a primit mesaje de care mesajul curent poate fi cauzal dependent, dar pe care P_d încă nu le-a primit;
- dacă $vt_m[s] > V_d[s] + 1$ atunci se întârzie m
 mai sunt mesaje de la P_s pe care P_s nu le-a primit
 (asigură ordinea FIFO pentru canale nonFIFO)
- dacă $vt_m[s] > V_d[s]$ atunci rejectează m
 m este un duplicat al unui mesaj primit anterior



Byzantine Generals

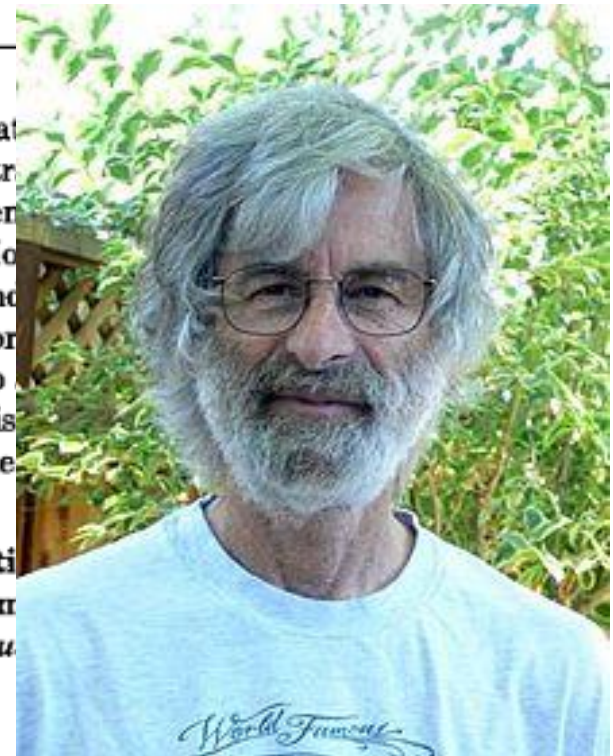
The Byzantine Generals Problem

LESLIE LAMPORT, ROBERT SHOSTAK, and MARSHALL PEASE
SRI International

Reliable computer systems must handle malfunctioning components that may fail arbitrarily and return inconsistent values to different parts of the system. This situation can be expressed abstractly as the Byzantine Generals Problem. Suppose n generals of the Byzantine army camped with their troops around an enemy city. A messenger can deliver a message from one general to another, but may be a traitor who will try to confuse the others. The problem is to find an algorithm by which the loyal generals will reach agreement. It is shown that, using only oral messages, the problem is solvable if and only if more than two-thirds of the generals are loyal; so for $n = 3$, two loyal generals. With unforgeable written messages, the problem is solvable for any number of loyal generals and possible traitors. Applications of the solutions to reliable computer systems are discussed.

Categories and Subject Descriptors: C.2.4. [Computer-Communications Systems]—*network operating systems*; D.4.4 [Operating Systems]: Communications—*network communication*; D.4.5 [Operating Systems]: Reliability—*fault tolerance*

General Terms: Algorithms, Reliability





Tipuri de defecte în sistemele distribuite

- Valoarea comună poate *să nu fie atinsă* din diferite cauze
 - pierderea mesajelor în mediul de comunicație
 - procesele pot produce rezultate greșite

- Tipuri de defecte ale proceselor
 - ***crash***: procesul devine nefuncțional
 - ***byzantine***: procesul trimite mesaje cu un conținut arbitrar

- O analogie pentru astfel de situații este *Problema Generalilor Bizantini*



Problema generalilor bizantini

- Câteva divizii ale *Armatei bizantine*, fiecare sub comanda unui *General*, înconjoară inamicul
- Generalii bizantini trebuie să ajungă la o înțelegere în privința atacului
- Este crucial să se ajungă la o înțelegere, deoarece doar un atac simultan din partea tuturor diviziilor poate conduce la o victorie
- Diviziile sunt dispersate geografic, așa că Generalii își comunică între ei observațiile asupra inamicului prin intermediul mesagerilor
- Astfel, cauzele care pot îngreuna stabilirea unei înțelegeri pot fi:
 - Mesagerii pot fi prinși de către inamic – așadar, mesajele nu ajung la destinație
 - Generalii pot fi **trădători**, încercând să deruteze Generalii **loiali**
- În continuare, vom presupune că mesagerii ajung la destinație



Problema generalilor bizantini

Generalii au nevoie de un algoritm care să garanteze că:

A: Toți generalii loiali decid același plan de acțiune.

Nu e suficient: vrem același plan pentru toți, însă dorim ca și planul să fie unul adecvat.

B: Un număr mic de trădători nu pot influența decisiv Generalii loiali.

Problema generalilor bizantini

- Fie $v(i)$ – informația comunicată de către Generalul cu numărul i
- Fiecare General folosește o metodă pentru a combina valorile $v(1), v(2), \dots, v(n)$ într-un plan de acțiune
- Condiția A este îndeplinită prin faptul că toți generalii folosesc aceeași metodă de a combina cele n valori
- Condiția B este îndeplinită prin asigurarea că metoda folosită este una robustă
- Dacă decizia care poate fi luată este *atac* sau *retragere*, $v(i)$ este opțiunea Generalului i dintre cele două variante, iar decizia finală pentru fiecare General se bazează pe votul majoritar (dintre cele n valori)



Problema generalilor bizantini

- Generalii își comunică unii altora valorile $v(i)$
- Generalii trădători pot trimite valori diferite celorlalți Generali
- Satisfacerea condiției A presupune ca fiecare General să ia decizia finală bazată pe aceeași mulțime $v(1), v(2), \dots, v(n)$
- Astfel, pentru îndeplinirea condiției A :
 1. Fiecare General loial trebuie să obțină aceleași informații $v(1), \dots, v(n)$.



Problema generalilor bizantini

1. Fiecare General loial trebuie să obțină aceleași informații $v(1), \dots, v(n)$.

- Condiția 1 implică faptul că un General poate să nu folosească valoarea $v(i)$ obținută direct de la Generalul i , de vreme ce un General trădător poate să trimită valori diferite
- Astfel, se introduce riscul ca un General să nu folosească o valoare $v(i)$, deși Generalul i este unul loial
- Se introduce următoarea cerință:

2. Dacă Generalul i este loial, atunci valoarea trimisă de el trebuie să fie folosită de către fiecare General loial.

- Condiția 1 se poate rescrie astfel:

1'. Oricare doi Generali loiali folosesc aceeași valoare pentru $v(i)$.



Problema generalilor bizantini

- Restrângem abordarea problemei la modul în care *un* General își trimite valoarea celorlalți Generali
- Termenii problemei se schimbă: un General comandant trimite ordinul Locotenenților, obținându-se astfel problema:
- **Problema Generalilor Bizantini:**
 - Un General comandant trebuie să trimită un ordin celor $n - 1$ Locotenenți astfel încât:
 - IC_1 : Toți Locotenenții loiali se supun aceluiași ordin
 - IC_2 : Dacă Generalul comandant este loial, atunci fiecare Locotenent loial se supune ordinului Generalului comandant
- IC_1, IC_2 – *interactive consistency conditions*
- Se observă că, pentru un Comandat loial, IC_1 rezultă din IC_2 ; totuși, Comandantul nu este în mod necesar loial

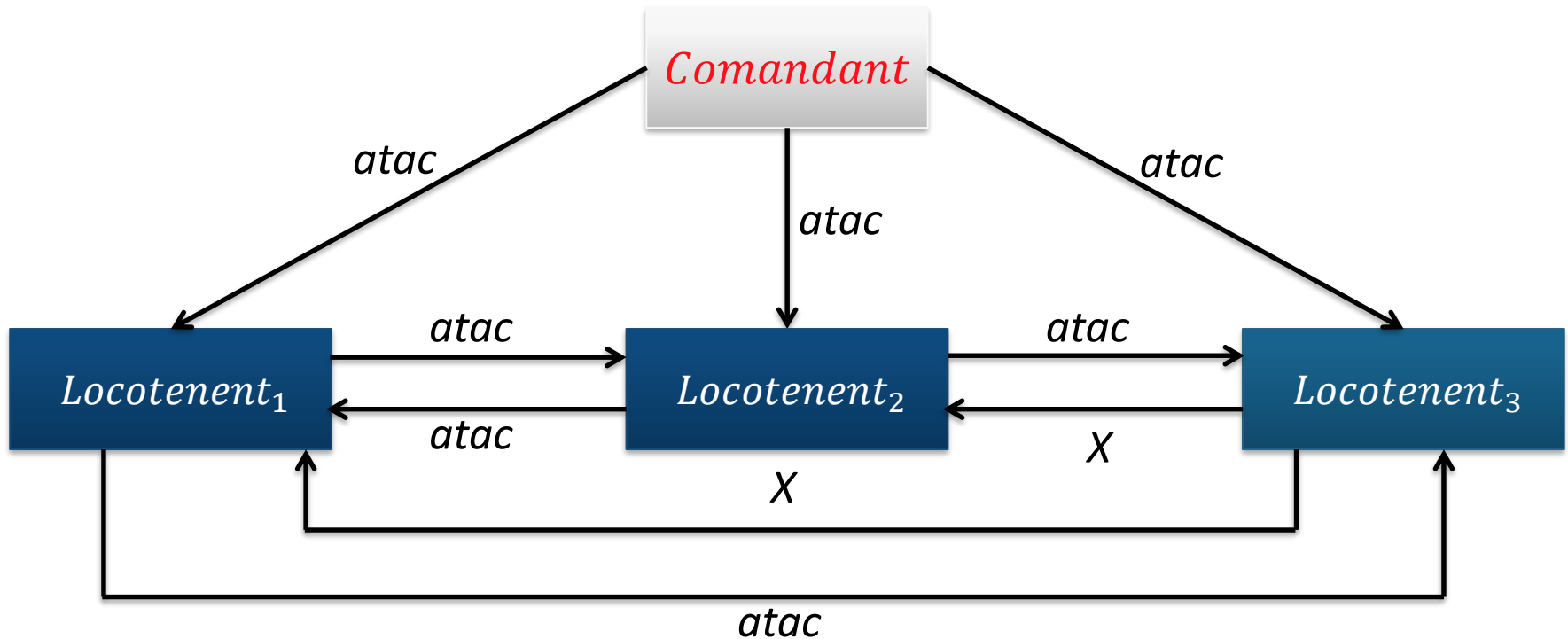


Problema generalilor bizantini

- Evenimentele sunt următoarele:
 - Comandantul trimite ordinul tuturor Locotenenților
 - Un Locotenent trimite celorlalți Locotenenți mesajul primit de la Comandant
 - După primirea mesajului de la Comandant și a copiilor de la ceilalți Locotenenți, un locotenent decide *prin vot majoritar* ce decizie va lua

Problema generalilor bizantini

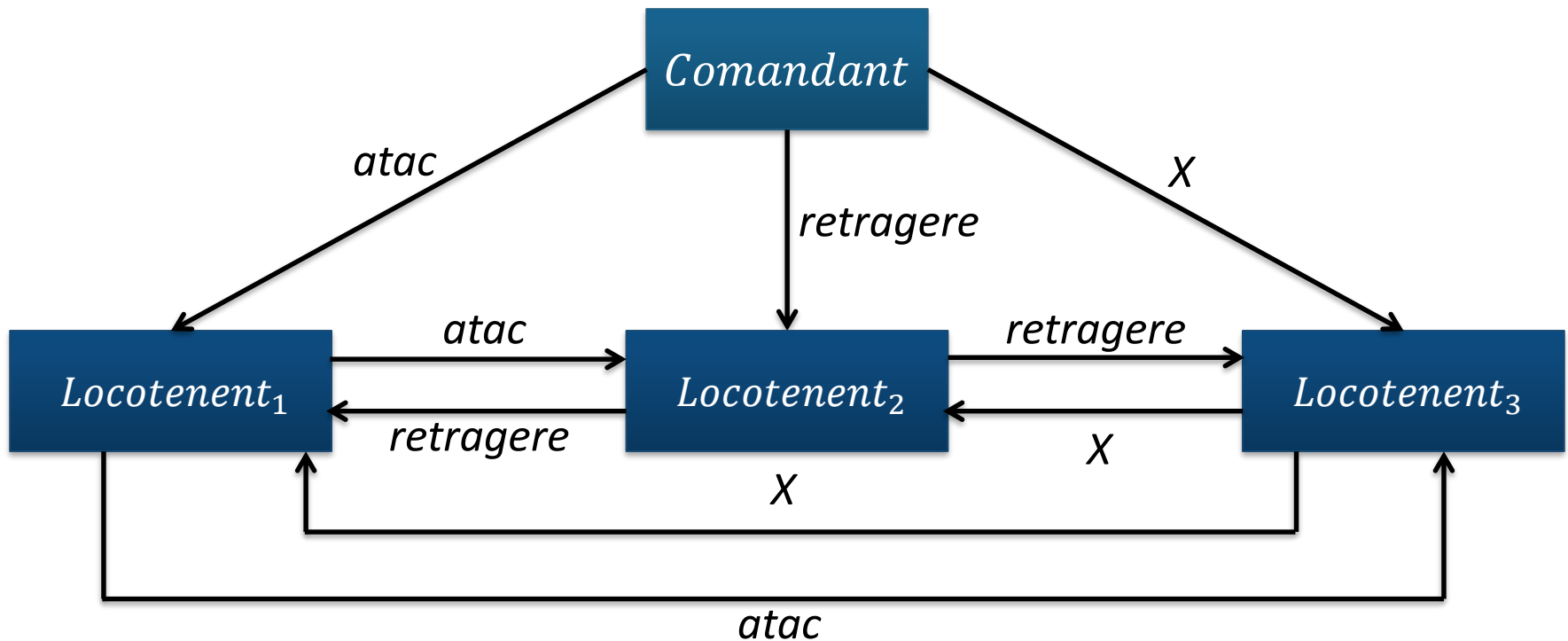
Exemplul 1



- **Locotenent₃** este trădător
- $X \in \{\text{atac}, \text{retragere}\}$
- Oricare ar fi mesajul transmis de trădător, cei doi Locotenenți loiali vor lua aceeași decizie (*atac*)

Problema generalilor bizantini

Exemplul 2



- *Comandant* este trădător
- $X \in \{\textit{atac}, \textit{retragere}\}$
- Oricare ar fi mesajul transmis de trădător, toți Locotenenții vor lua aceeași decizie X

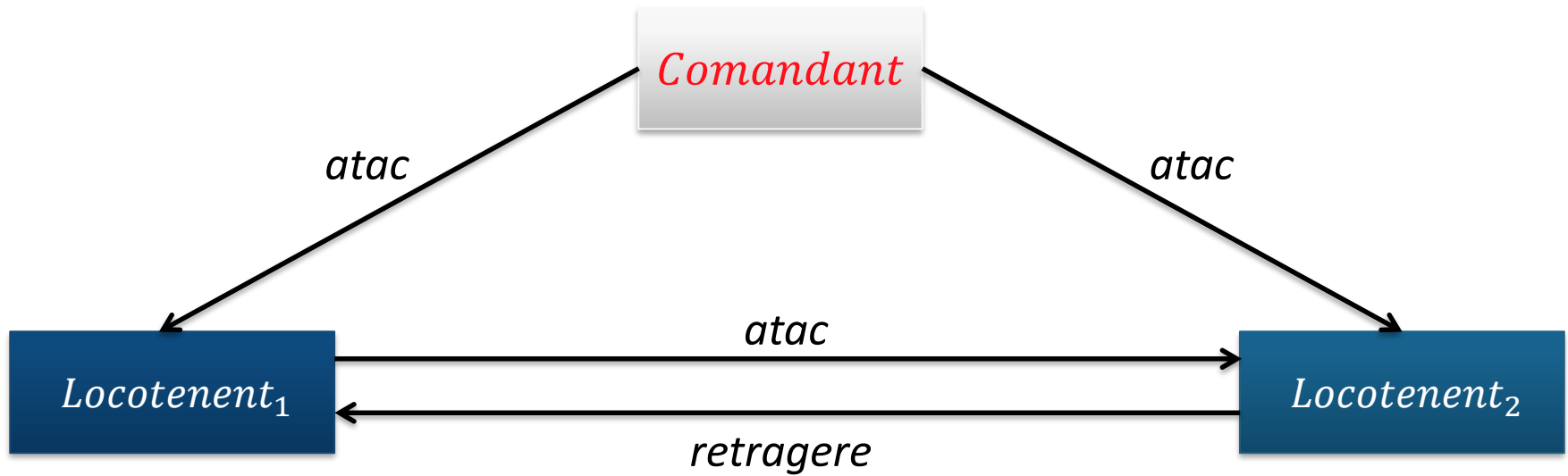


Problema generalilor bizantini

- Dificultatea problemei constă în faptul că:
 - Dacă Generalii (Locotenenții) pot trimite doar mesaje orale, atunci nu există soluție decât pentru cazul în care *2/3 din Generali sunt loiali*
- Un mesaj *oral* este aflat complet sub controlul emițătorului, deci un trădător poate trimite orice mesaj
- Mesajul trimis este *atac / retragere*
- Următoarele exemple: nu există soluție pentru 3 Generali, din care un trădător

Problema generalilor bizantini

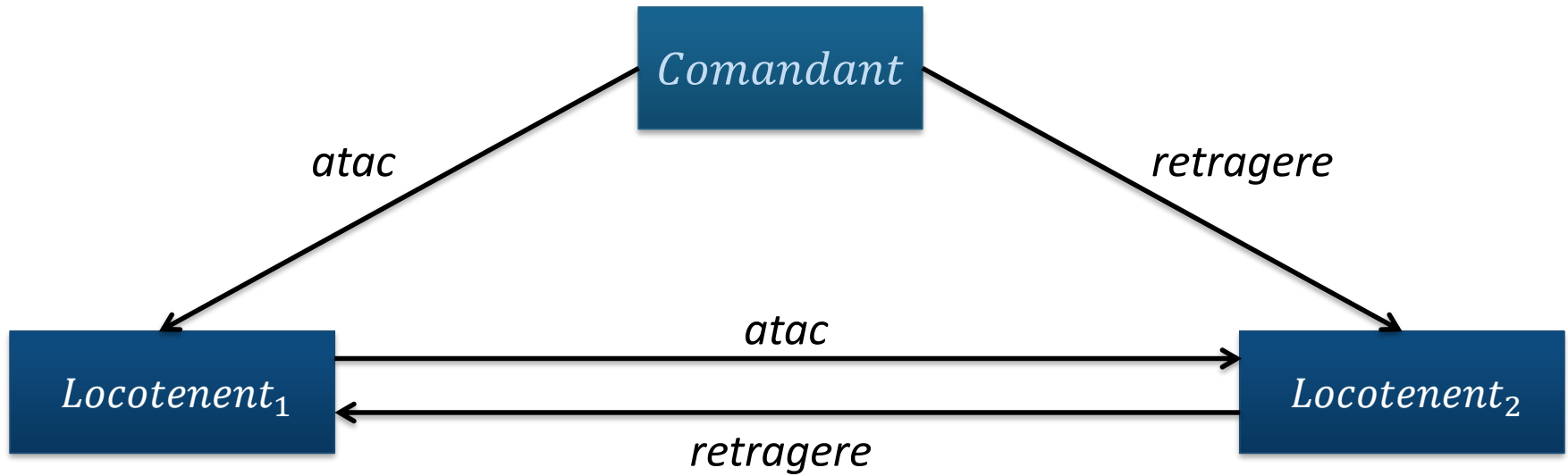
Exemplul 3



- *Locotenent₂* este trădător

Problema generalilor bizantini

Exemplul 4



- *Comandant* este trădător

Problema generalilor bizantini

Soluția cu mesaje orale

- Mesajele orale îndeplinesc următoarele condiții:
 - A_1 : Fiecare mesaj trimis ajunge corect la destinație
 - A_2 : Receptorul mesajului cunoaște autorul mesajului
 - A_3 : Absența unui mesaj poate fi detectată
- Fiecare General (Locotenent) poate trimite mesaje oricărui alt General
- $A_1, A_2 \rightarrow$ un trădător nu intervine în comunicarea dintre alți doi Generali
- $A_3 \rightarrow$ un trădător nu poate influența decizia prin a nu trimite un mesaj
- Un Comandant trădător poate decide să nu trimită ordine, așa că este nevoie de o decizie implicită pentru Locotenenți
- Decizia implicită este: *retragere*

Problema generalilor bizantini

Soluția cu mesaje orale

- Se definește algoritmul pentru soluția cu mesaje orale: $OM(m)$
- $OM(m)$ rezolvă Problema Generalilor Bizantini pentru minim $3m + 1$ Generali și cel mult m trădători
- Algoritmul folosește funcția $majority(v(1), \dots, v(n - 1))$
 - Dacă majoritatea valorilor $v(i)$ este v , atunci $majority$ este v
- Funcția $majority(v(1), \dots, v(n - 1))$ va returna:
 - Valoarea majoritară din $v(i)$, dacă aceasta există; altfel, se returnează *retragere*
 - Mediana valorilor $v(i)$, presupunând că mulțimea este sortată



Problema generalilor bizantini

Soluția cu mesaje orale

■ $OM(0)$

- 1) Comandantul își trimite valoarea tuturor Locotenenților.
- 2) Fiecare Locotenent folosește valoarea primită de la Comandant, sau *retragere* dacă nu primește nimic.

■ $OM(m), m > 0$

- 1) Comandantul își trimite valoarea tuturor Locotenenților.
- 2) Pentru fiecare i , fie $v(i)$ valoarea primită de Locotenentul i de la Comandant (sau valoarea implicită – *retragere*). Locotenentul i ia rolul de Comandant în $OM(m - 1)$, pentru a trimite valoarea celorlalți Locotenenți.
- 3) Pentru fiecare i , pentru fiecare $j, j \neq i$, fie $v(j)$ valoarea pe care Locotenentul i a primit-o de la Locotenentul j în pasul 2) ($OM(m - 1)$) sau valoarea implicită.
Locotenentul i folosește *majority*($v(1), \dots, v(n - 1)$).

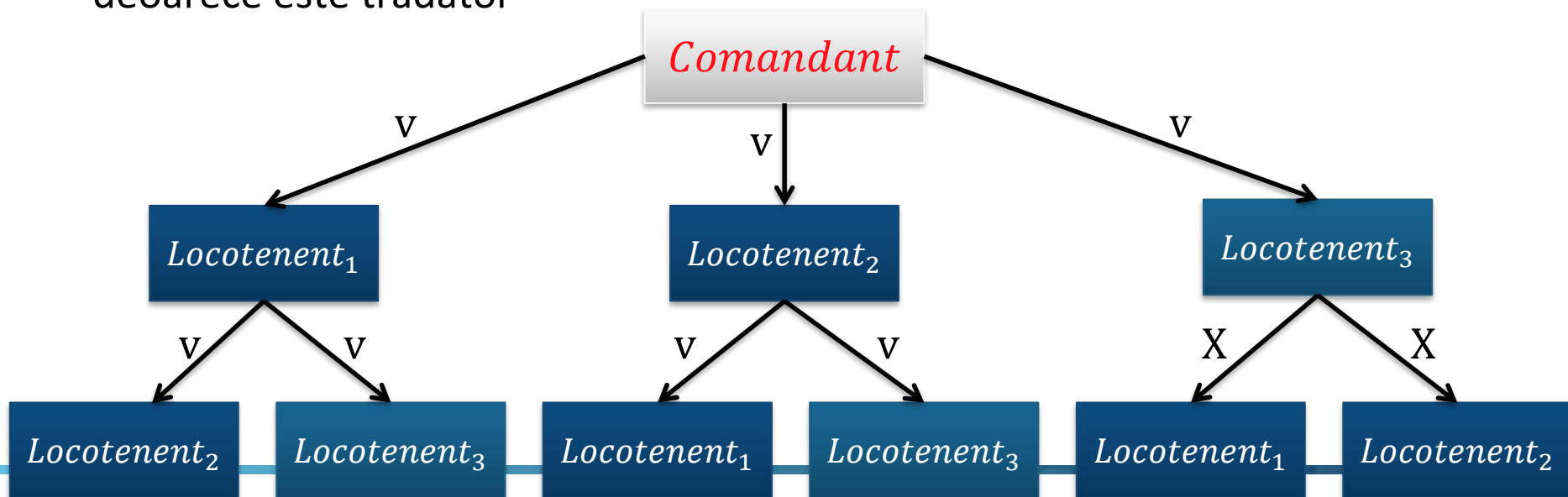


Soluția cu mesaje orale Exemplu

- $m = 1, n = 4 \rightarrow 1$ Comandant, 3 Locotenenți (Locotenentul cu nr.3 = trădător)
- $OM(1)$
 - Comandantul trimite valoarea v
 - Pentru fiecare i
 - $v(1) = v$ $Locotenent_1 = \{v\}$
 - $OM(0)$ $Comandant = Locotenent_1$
 - Comandantul își trimite valoarea v
 - $Locotenent_2 = \{v\}$
 - $v(2) = v$ $Locotenent_2 = \{v, v\}$
 - $OM(0)$ $Comandant = Locotenent_2$
 - Comandantul își trimite valoarea v
 - $Locotenent_1 = \{v, v\}$
 - $v(3) = v$
 - $OM(0)$ $Comandant = Locotenent_3$
 - Comandantul își trimite valoarea X
 - $Locotenent_1 = \{v, v, X\}$
 - $Locotenent_2 = \{v, v, X\}$

Soluția cu mesaje orale Exemplu

- Pentru fiecare i
 - $Locotenent_1$ folosește $majority(v, v, X) = v$
 - $Locotenent_2$ folosește $majority(v, v, X) = v$
 - $Locotenent_3$ este trădător
- Observație: parcursul pentru cel de-al treilea Locotenent nu a fost urmărit, deoarece este trădător



Soluția cu mesaje orale Complexitate

■ Complexitate:

- Pas 1: $(n - 1) \times OM(m - 1)$ mesaje
- Pas 2: $(n - 1) \times (n - 2) \times OM(m - 2)$ mesaje
- Pas 3: $(n - 1) \times (n - 2) \times (n - 3) \times OM(m - 3)$ mesaje
- Deci, în pasul k : $(n - 1) \times (n - 2) \times \dots \times (n - k) \times OM(m - k)$ mesaje
 - $k = m$: $(n - 1) \times (n - 2) \times \dots \times (n - m) \times OM(0)$
 - Pasul $m + 1$: se trimit $(n - 1) \times (n - 2) \times \dots \times (n - m - 1)$ mesaje
 - Număr total de mesaje: $O(n^{m+1})$



Corectitudine

LEMA 1. Pentru orice m și k , $UM(n, m)$ satisface **IC2** dacă numărul n de generali este mai mare de $2k+m$ și sunt cel mult k trădători.

TEOREMA 1. Pentru orice m , algoritmul $UM(n, m)$ satisface condițiile IC1 și IC2 dacă numărul de generali n este mai mare de $3m$ și sunt cel mult m trădători.

IC1. Toți locotenenții loiali se supun aceluiași ordin.

IC2. Dacă comandantul este loial, atunci fiecare locotenent loial se supune ordinului transmis de acesta.

ipoteză: $n > 2k + m \rightarrow n - 1 > 2k + (m - 1)$

(ipoteza inducție satisfăcută \rightarrow fiecare locotenent loial obține $v_j = v$ pentru fiecare locotenent loial j .

Sunt cel mult k trădători și $n - 1 > 2k + (m - 1) \geq 2k$
 \rightarrow majoritatea dintre $n-1$ locotenenți sunt loiali

ndant, sau folosește V_{def}

Algoritmul $UM(n, m)$, $m > 0$.

(1) Comandantul trimite valoarea sa fiecărui

\rightarrow fiecare locotenent loial obține o majoritate de valori $v_j = v \rightarrow$ în pasul (3) se obține $majority(v_1, \dots, v_{n-1}) = v \rightarrow IC2$ satisfăcută

Lema se referă la $IC2$ (comandant loial). Dem inductie. Pentru $m=0$: comandant loial trimite v ; fiecare mesaj transmis este livrat corect (cf. propr. A1) $\rightarrow UM(n, 0)$ satisface $IC2$

ceilați $n - 2$ locotenenți folosind $UM(n - 1, m - 1)$

(2) For each i and each $j > i$:

În pas (1), comandant loial trimite v celor $n-1$ locotenenți.

În (2), fiecare locotenent loial aplică $UM(n-1, m-1)$ cu $n-1$ generali.

Pp. proprietatea $IC2$ indeplinită ptr. $m-1$, $m > 0$ și probăm ptr m .

atul i o primește de la Locotenentul j în pasul (2) $(v_j, m - 1)$), sau V_{def} dacă nu primește nici o

a $majority(v_1, \dots, v_{n-1})$.



Algoritmul UM(n,0).

TEOREMA 1. Pentru orice m , algoritmul UM(n,m) satisface condițiile IC1 și IC2 dacă numărul de generali n este mai mare de $3m$ și sunt cel mult m trădători.

Presupunem teorema adevărată ptr $m-1$; probăm pentru $m > 0$.

Comandant loial: punem $k=m$ în Lema 1 ($n > 2k+m \sim n > 3m$) \rightarrow UM(n,m) satisface IC2, iar IC2 \rightarrow IC1

locotenent loial se supune ordinului transmis de el).

(2) **For each** Locotenent i

fie v_i valoarea primită de la comandant. sau V_{def} dacă nu primește nici o valoare.

Dem IC1 ptr. comandant trădător:

Max. m trădători și comandant trădător \rightarrow max. $m-1$ locotenenți trădători.

Sunt $> 3m$ generali $\rightarrow > 3m-1$ locotenenți. Deoarece $3m-1 > 3(m-1) \rightarrow$ ipoteza de inducție satisfăcută \rightarrow UM(n-1,m-1) satisface condițiile IC1 și IC2.

\rightarrow ptr. fiecare j , orice doi locotenenți loiali obțin aceeași valoare v_j în pasul (3).

\rightarrow obțin aceleași valori $v_1 \dots v_{n-1} \rightarrow$ calculează aceeași valoare $\text{majority}(v_1 \dots v_{n-1})$ în pasul (3)



Soluția cu mesaje semnate

- Se adaugă condiția:
 - A_4
 - a) Semnătura unui general loial nu poate fi falsificată și orice alterare a mesajelor sale semnate poate fi detectată.
 - b) Oricine poate verifica autenticitatea unei semnături.
- Trădătorii își pot falsifica semnăturile între ei
- Algoritmul cu mesaje semnate rezolvă problema pentru orice număr de Generali

Soluția cu mesaje semnate

- Desfășurare algoritm:
 - Comandatul trimite un mesaj semnat tuturor Locotenenților
 - Fiecare Locotenent își adaugă semnătura și trimite mesajul celorlalți Locotenenți

- Utilitatea semnăturii (de exemplu):
 - Comandatul este trădător și trimite mesaje diferite Locotenenților
 - Locotenenții, prin intermediul semnăturii, pot observa că mesajele diferite primite sunt cauzate chiar de Comandant – două mesaje diferite conțin semnătura Comandantului



Soluția cu mesaje semnate

- Algoritmul folosește o funcție *choice* care este aplicată unei mulțimi de ordine – pentru a se obține un *singur ordin*
- Cerințe:
 - Dacă mulțimea V conține un singur ordin v , atunci $choice(V) = v$
 - $choice(\emptyset) = RETRAGERE$, \emptyset este mulțimea vidă
- O posibilă definiție pentru *choice* este elementul median al mulțimii V – presupunând că mulțimea este ordonată



Soluția cu mesaje semnate

■ Notatii:

- $x : i$ – valoarea x semnată de Generalul i
 - $v : j : i$ – valoarea v semnată de Generalul j , apoi de Generalul i
- Generalul 0 este Comandantul
- V_i – mulțimea de ordine primite (corect semnate) de către Generalul i
 - Pentru un Comandant loial, această mulțime conține un singur element (Locotenenții loiali pot recunoaște mesajele false introduse de trădătorii)

- A nu se confunda *mulțimea de ordine* V_i cu mulțimea de mesaje primite de un General (mai multe mesaje pot conține același ordin)



Soluția cu mesaje semnate

Algoritmul $SM(m)$

- $V_i = \emptyset$
- Comandatul își semnează valoarea și o trimite fiecărui Locotenent
- Pentru fiecare i :
 - Dacă Locotenentul i primește un mesaj de forma $v : 0$ de la Comandant și mulțimea sa de ordine este vidă:
 - $V_i = v$
 - trimite $v : 0$: i celorlalți Locotenenți
 - Dacă Locotenentul i primește un mesaj de forma $v : 0 : j_1 : \dots : j_k$ și v nu este în mulțimea V_i , atunci:
 - se adaugă v la V_i
 - dacă $k < m$, se trimite mesajul $v : j_1 : \dots : j_k : i$ Locotenenților care nu sunt în mulțimea $\{j_1, j_2, \dots, j_k\}$
- Pentru fiecare i : atunci când Locotenentul i nu mai primește mesaje, se va supune ordinului $choice(V_i)$



Soluția cu mesaje semnate

Algoritmul $SM(m)$

- Cum se detectează faptul că un Locotenent nu va mai primi mesaje?
 - folosirea unui timer
 - se poate arăta că pentru o secvență $\{j_1, j_2, \dots, j_k\}$, $k \leq m$, un Locotenent poate primi cel mult un mesaj de forma $v : 0 : j_1 : \dots : j_k$
 - se poate impune ca Locotenentul j_k să trimită un astfel de mesaj sau un mesaj care să indice că nu va trimite un astfel de conținut

Corectitudine

TEOREMA 2. Pentru orice m , Algoritmul $SM(m)$ rezolvă problema generalilor bizantini dacă există cel mult m trădători.



Dacă comandantul este **loial**, atunci fiecare locotenent loial se supune ordinului transmis de el.

Inițial $V_i = \Phi$.

(1) Comandantul semnează și trimite valoarea sa fiecărui locotenent

(2) **For each** i :

Pas 1. Comandantul trimite ordinul semnat $v:0$

(A) **If** Locotenent i primește un mesaj de forma v :
primit încă nici un ordin **then**

(i) $V_i := \{v\}$;

(ii) transmite mesajul $v:0:i$

Fiecare locotenent loial primește ordinul v în pasul (2)(A).

(B) **If** Locotenent i primește un

(i) adaugă v la V_i ;

(ii) **if** $k < m$ **then** trimite mesaj $v:0:j_1: \dots :j_k:i$ fiecărui locotenent diferit de $j_1 \dots j_k$.

(3) **For each** i :

when Locotener

Un locotenent neloial nu poate falsifica ordinul comandantului → un locotenent loial nu poate primi un alt ordin în pasul (2)(B)



1. Toți locotenenții loiali se supun aceluiași ordin. Analizăm cazul "comandant **tradator**".

Inițial $V_i = \Phi$.

(1) Comandantul semnează și trimite valoarea sa fiecărui locotenent.

(2) **For each** i :

(A) **If** Locotenent i primește un mesaj de forma v în pasul (2)(A)(ii); j îl primește (conform proprietății A1).

Dacă i adaugă v la V_i în pasul (2)(B)(i), atunci el trebuie să fi primit un mesaj de forma $v:0:j1: \dots :jk$. Dacă j este unul dintre j_r , atunci (cf A4) el trebuie să fi primit deja ordinul v .

Dacă j nu este unul din j_r & $k < m$: i trimite mesajul $v:0:j1: \dots :jk:i$ lui j
→ j trebuie să primească ordinul v .

(B) **If** Locotenent i primește un mesaj de forma $v:0:j1: \dots :jk$ și v nu este în V_i **then**

(i) adaugă v la V_i ;

(ii) **if** $k < m$ **then** trimite mesaj $v:0:j1: \dots :jk:i$ fiecărui locotenent diferit de $j1$

(3) Dacă j nu este unul din j_r & $k = m$: comandant trădător → cel mult $m - 1$ locotenenți sunt trădători. Cel puțin unul dintre $j1, \dots, jm$ este loial. El a trimis lui j valoarea v atunci când a primit-o prima dată.

→ dacă i are v atunci și j are v → toți lt. loiali au aceleași v -uri

Soluția cu mesaje semnate

Exemplu

