

Parameterised Light Atomic Broadcast and Membership Protocol (PLAMP)

Jean Fanchon

LAAS-CNRS

Introduction

Besoins

Problèmes

L'algorithme

Sur réseau synchrone fiable

Crash de processus, pertes de messages

Variables, paramétrage

Implémentation et architecture

Conclusions

Coopération, coordination d'applications distribuées

Tolérance aux fautes

- [**Membership**
 - Listes de membres, join quit
 - Notifications d'événements, nouveaux, processus fautifs
- [**Atomic Broadcast**
 - Diffusion de messages
 - Livraison des mêmes messages dans le même ordre (total)
- [**Utilisations**
 - Gestion de replica
 - Partage d'objets par des applications distribuées
 - Gestion distribuée de ressources
 - Allocation unique de noms, décisions, etc....

Ex: TSF OLC RST

Problèmes similaires mais contextes et utilisations différentes

--> **Code ouvert, homogène, accessible, adaptable**

Problème du consensus (AB, Memb, etc..)

Difficultés intrinsèques

Asynchrone: impossibilité avec 1 seul processus fautif (FLP)

Synchrone : t fautes requièrent t+1 phases (T)

Réseau non-fiable

---> Restreindre les hypothèses de fautes

---> Décomposition du problème et résolution par modules (FD)

Difficultés d'implémentation

Chaque module nécessite des échanges de messages

Crash de composants vs. crash de sites

Interfaces, codes hétérogènes, maintenance, taille..

Modules / Fonctions

Atomic broadcast

Group membership

Consensus

Reliable communication

Point to point communication

**Failure
Dectectors**

Transport Protocol

Membership

**Agreement on group membership: if two correct processes belong to group
then they have the same view of the membership**

**Total ordering of change notifications: any two processes deliver
membership changes notifications in the same order**

**Virtual synchrony : membership changes notifications are delivered
at the same point in the application message streams**

Atomic Broadcast

B-Termination : if a correct process executes **broadcast(m)**
then all correct process execute **delivery(m)**

B-Validity : if a process executes **delivery(m)**
then m has been sent by a process

RB-Termination : if a process executes **delivery(m)**
then all correct processes execute **delivery(m)**

ORB-Agreement: if two processes p and q deliver m1 and m2
p and q deliver m1 and m2 in the same order

Principes de PLAMP (1)

Broadcast ordre total avec réseau et processus fiables

chaque processus i | identificateur unique

| liste des membres

Phase = i , chaque processus émet vers tous les autres

Chaque processus reçoit de tous les autres

Phase = phase + 1

Algorithme

par phases

Pour tous processus k et h | $|\text{phase}(k) - \text{phase}(h)| \leq 1$

A la fin de chaque phase, livraison de la même liste de messages

Principes de PLAMP (2)

Prise en compte de crash de processus -->Failure detector P

**Crash de processus: arrête d'émettre en cours de phase
n'émet plus phases suivantes**

Réseau fiable synchrone

Algorithme

modifié

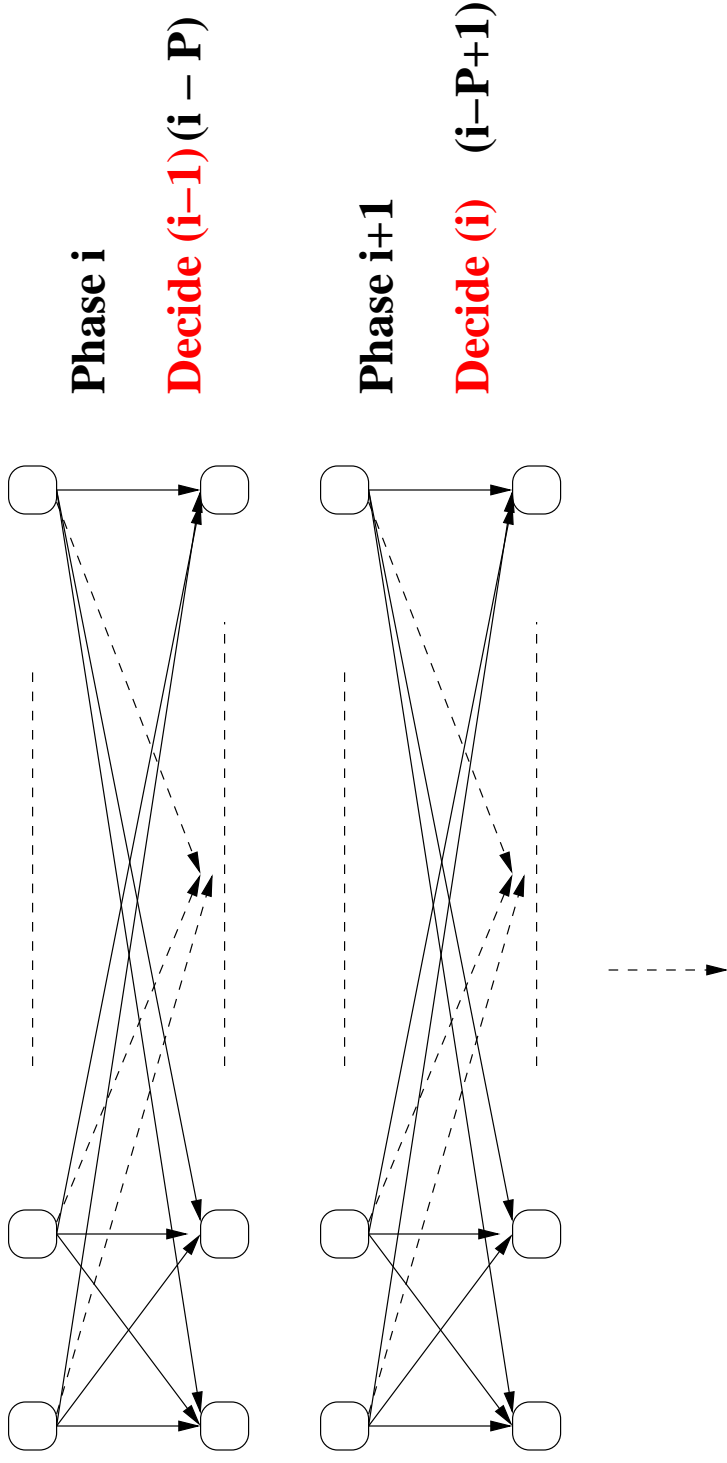
Fin de phase: réception de messages de tous ou bien timeout

Détection de processus fautifs

Message de la phase i acquittement de la phase $i-1$, liste de fautifs

A la fin de phase i, décision quant à la phase $i-1$ (i-P)

Phases et décisions



Decision: 1) processus fautifs 2) liste de messages à livrer

Variables locales

j,k,l,h: ident

n: phase

Rec_k(-)(n)

l	h
1	0

1 = received
0 = not received

Val_k(-,)(n)

1/0				
1/0	X			
1/0			X	

} idents qui ont
 recu de j
 en phase n-1

Val_k(-j)(n) j

Req_l(-)(n-1) Req_h(-)(n-1)

Décisions

Pas 2 phases consécutives avec faute(s)

Décision pour phase $i-1$

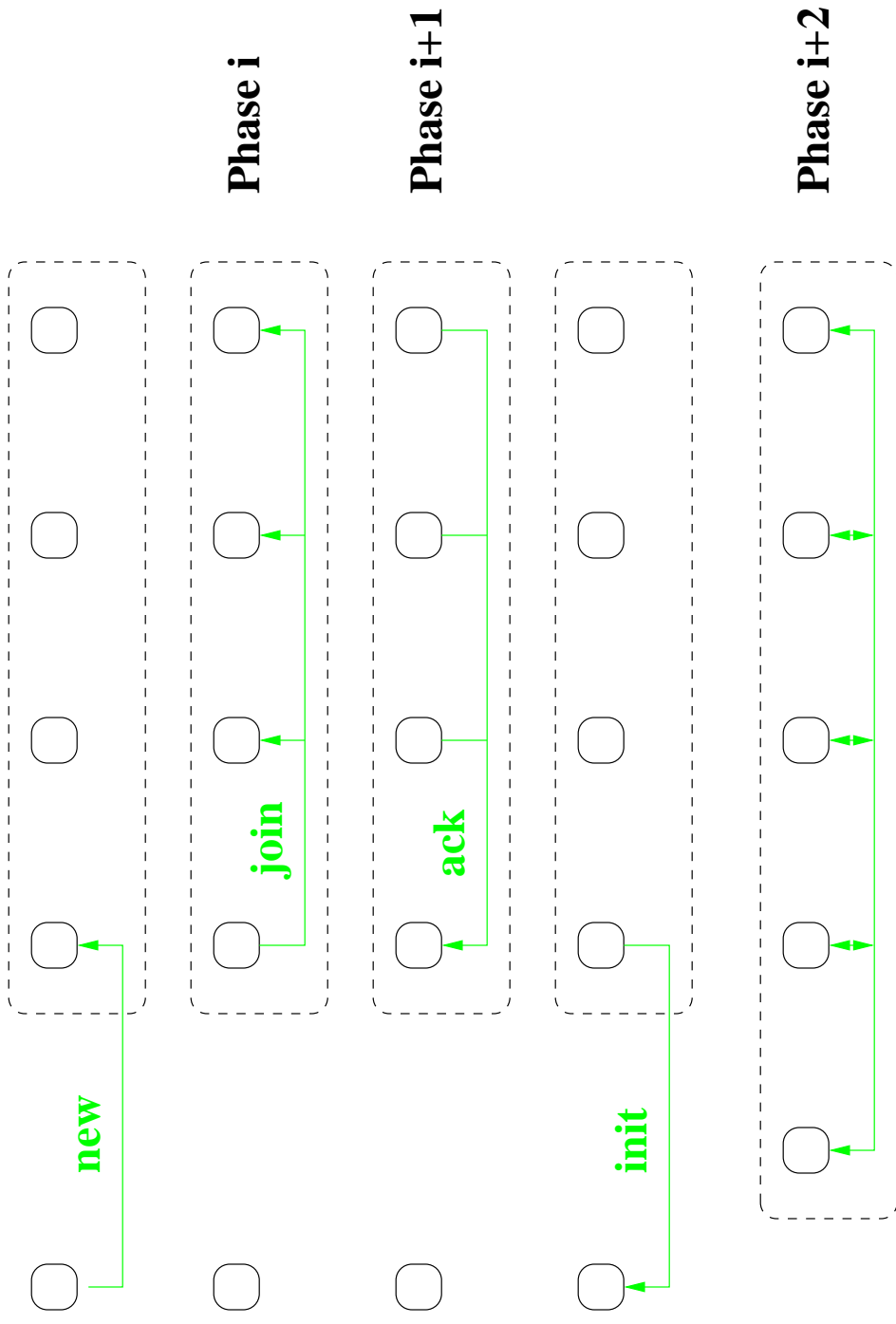
$$| \mathbf{h}, \mathbf{Val}_k(\mathbf{h}, \mathbf{j})(\mathbf{n}) = \mathbf{0} | = \mathbf{0}$$

Décision sur phases successives dissymétriques

If $f < n/2$ | $\{ \mathbf{h}, \mathbf{Val}_k(\mathbf{h}, \mathbf{j})(\mathbf{n}) = \mathbf{1} \} | > n/2 \Rightarrow$

For all l | $\{ \mathbf{h}, \mathbf{Val}_l(\mathbf{h}, \mathbf{j})(\mathbf{n}) \in \{ \mathbf{0}, \mathbf{X} \} \} | < n/2$

Join scenario



Principes de PLAMP (3)

Réseau non fiable, asynchrone **--> FD** **◇ S**

En cours de phase, demandes de réémissions
réémissions de message phase i ou i-1

Fin de phase: réception de messages de tous ou bien timeout

Détection de processus fautifs ou de messages perdus

Message de la phase i acquittement de la phase i-1

A la fin de phase i , décision quant à la phase i-1 (i-P)

Algorithme

modifié

Décisions

Pas 2 phases consécutives avec faute(s)

Décision pour phase $i-1$

$$|h, \text{Val}_k(h,j)(n) = 0 \mid \neq 0 \text{ ET } \text{Rec}_k(j)(n) = 0$$

Pas t phases consécutives avec faute(s)

Décision pour phases $i-1$ à $i-t$

Décision sur phases successives dissymétriques

$$\text{If } f < n/2 \quad | \{ h, \text{Val}_k(h,j)(n) = 1 \} | > n/2 \Rightarrow$$

$$\text{For all } l \quad | \{ h, \text{Val}_l(h,j)(n) \in \{ 0, X \} \} | < n/2$$

Paramétrage

Timeout et nombre des demandes de réémissions

Décision du failure detector

Nombre de phases utilisées

Protocole de transport (UDP → TCP)

→ Adaptable au réseau

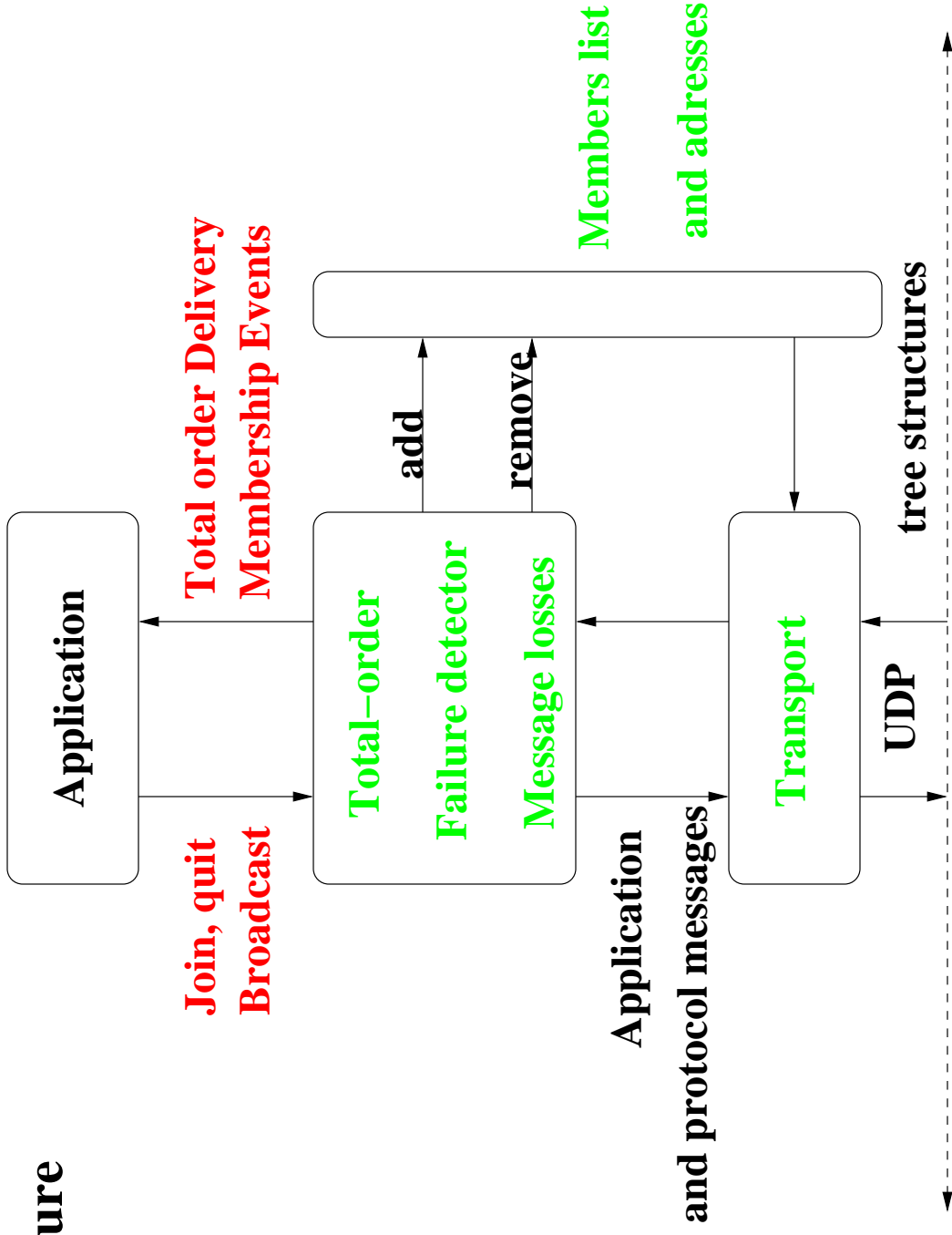
au type d'applications

au degré de développement

→ Paramétrage à la main (FD)

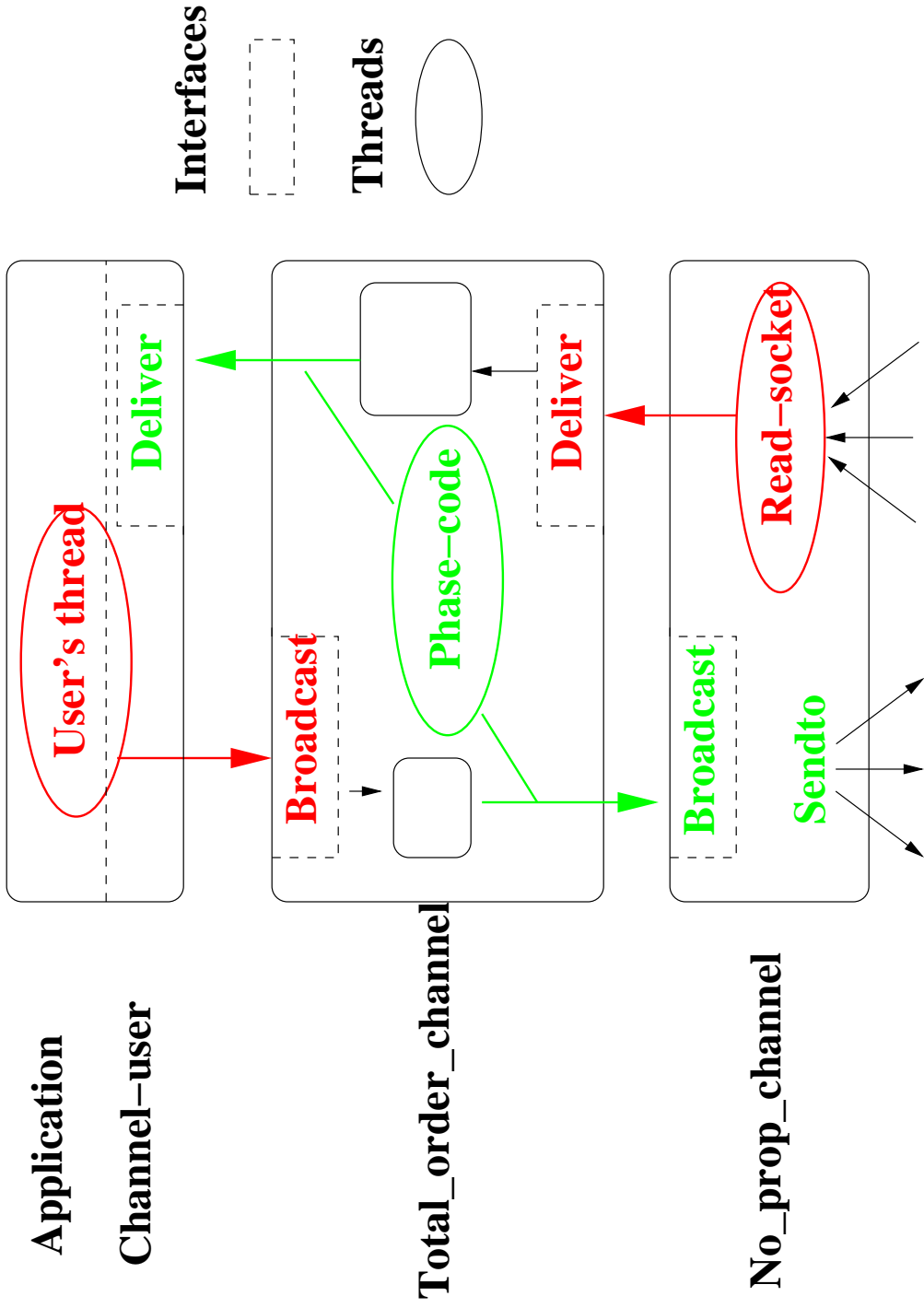
autostabilisation (timeouts)

Structure



TCP, IP multicast, etc....

Classes et threads



Conclusions

Membership et atomic broadcast sur réseau non fiable

Paramétrable , ouvert, compact

Taille code: 1500 lignes C++

Taille exécutable : 40KB

Outil de test pour algorithmes

Crash et recovery:implémentable au niveau applicatif

Compatible avec couche d' authentification

Arborescence de diffusion.

Failure Detectors

Completeness Tout processus crashé est suspecté par tous les corrects

Accuracy

◇ **S** Tout processus correct n'est plus suspecté après un certain temps

P^x Un processus ne suspecte faussement que au plus $n-x-1$ corrects

P Un processus n'est jamais suspecté avant de crasher

P <----- Interactive consistency $D_k[i] \in \{v_i, \perp\}$

◇ **S** <----- Consensus pour $t < n/2$

