

Exclusion mutuelle distribuée, algorithmes arborescents et protocoles Ethernels

Gérard Le Lann
INRIA – IMARA
Gerard.Le_Lann@inria.fr

Je n'aime pas les gens qui parlent pendant que je les interromps ...

Exclusion mutuelle **distribuée non bloquante** : hypothèses

Soit un jeu à n participants, noms uniques $\in [1, n]$, $n > 1$, qui doivent utiliser un objet partagé C . Un joueur dispose de deux actions sur C :

- ❖ A , activée pour tenter d'acquérir C (la durée d'utilisation de C est finie bornée)
- ❖ L , pour libérer C après utilisation.

C peut prendre **3** états, 0 (libre), 1 (utilisé), ∇ (**conflit**), notés $e(C)$.

Une action A exécutée par j à t produit l'une des réponses suivantes :

□ **Gagné** (C est alloué à j) :

- $e(C) = 0$ à t , j est le seul demandeur || $e(C) : 0 \Rightarrow 1$

□ **Perdu** (C n'est pas alloué à j) :

- $e(C) = 1$ à t || $e(C)$ reste à 1 (jusqu'à L par $k \neq j$, et alors $e(C) \Leftarrow 0$)

□ **Perdu** (C n'est pas alloué à j) :

- $e(C) = 0$ à t , j n'est pas le seul demandeur à t || $e(C) : 0 \Rightarrow \nabla$ || $e(C) = \nabla$ pendant une durée σ , après quoi $e(C) \Leftarrow 0$

- $e(C) = \nabla$ à t || $e(C)$ reste à ∇

Les joueurs observent les transitions d'état et les états de C , mais ne peuvent se concerter. A activable librement. **Tout joueur peut s'arrêter à tout moment.**

Exclusion mutuelle distribuée **non bloquante** : propriétés

- ❖ Sûreté logique : à tout instant, au plus 1 joueur utilise C
- ❖ Vivacité : tout joueur qui active A se voit allouer C en un temps fini*
- ❖ Non oisiveté : $e(C) \neq 0$ si au moins 1 joueur a activé A et C ne lui est pas encore alloué.

*Un joueur qui utilise C peut ne pas exécuter L en un temps fini (arrêt accidentel).

Problème classique (Dijkstra, ..., 60's) en version centralisée : exclusion mutuelle pour entrer en section critique (programmes concurrents).

Solution (*sucre syntaxique*) : variable *mutex*, virtualisation de C, $e(mutex) = \{0, 1\}$

Depuis les 70's, nombreuses déclinaisons du problème (et des solutions) en version distribuée :

- Sémaphores
- Estampilles
- **Jeton circulant**
- **Tris arborescents**

© IFIP, NORTH-HOLLAND PUBLISHING COMPANY
(1977)

<https://www.rocq.inria.fr/novaltis/.../IFIP%20Congress%201977.pdf>

Format de fichier: PDF/Adobe Acrobat - [Afficher](#)
de G LE LANN - 1977 - [Cité 428 fois](#) - [Autres articles](#)
GERARD LE LANN. IRISA—Université dc ... problems is
illustrated by the study of a *mutual exclusion* scheme
intended for a *distributed* envi- ronment. 1 .

Tris arborescents // Tree searches

Buts : en cas de conflit, trier les joueurs permet de leur affecter des instants d'activation à nouveau de A qui sont différents (élimination de la concomitance)

Principe : isoler des sous-arbres, chacun d'entre eux contenant au plus 1 joueur actif

Disciplines : théorie de l'information, structures de données, algorithmique distribuée, modélisation/calcul analytique, analyse combinatoire, ...

• Algorithmes probabilistes

Les joueurs choisissent leurs parcours selon une même loi aléatoire → les parcours peuvent ne pas différer (probabilité calculable) → probabilité non nulle de violation de la propriété de vivacité (a fortiori, pas de borne de terminaison)

• Algorithmes déterministes

Les joueurs choisissent leurs parcours selon une même loi déterministe → les parcours finissent sûrement par différer → terminaison certaine en temps fini borné

• Analyses comportementales // performances (débits, délais)

Stochastiques : moments de distribution (espérances, écart-types, ...), comportements asymptotiques

Pires cas : bornes inf (débits), bornes sup (délais)

```

FUNCTION GetQuorum (Tree: NetworkHierarchy): QuorumSet;
  VAR left, right : QuorumSet;
  BEGIN
  IF Empty (Tree) THEN
    RETURN ({});
  ELSE IF GrantsPermission(Tree↑.Node) THEN
    RETURN ((Tree↑.Node) ∪ GetQuorum (Tree↑.LeftChild));
    OR
    RETURN ((Tree↑.Node) ∪ GetQuorum (Tree↑.RightChild));(*line 9*)
  ELSE
    left←GetQuorum(Tree↑.left);
    right←GetQuorum(Tree↑.right);
    IF (left = ∅ ∨ right = ∅) THEN
      (* Unsuccessful in establishing a quorum *)
      EXIT(-1);
    ELSE
      RETURN (left ∪ right);
    END; (* IF *)
  END; (* IF *)
  END; (* IF *)
END GetQuorum

```

Figure: Algorithm for constructing a tree-structured quorum.

Preuves de propriétés (cf. # 3) requises pour toute solution/algorithmme

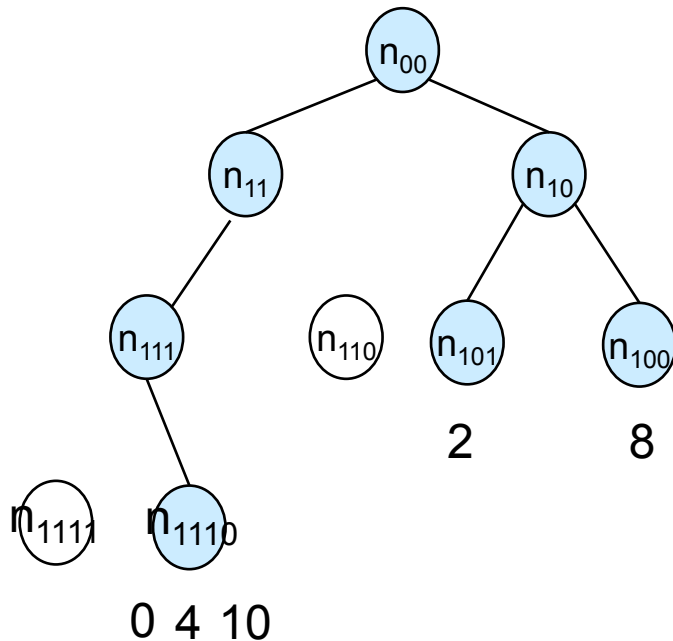
(« home work »)

Exemple du tri binaire – Algorithme probabiliste

A chaque réponse « perdu », tirage aléatoire sur $\{0, 1\}$

0 → tenter A à nouveau à $t = \text{transition } e(C) \leftarrow 0$

1 → tenter A à nouveau à $t + \sigma$ et attendre transition $e(C) \leftarrow 0$



1, n_{00} : $e(C) = 0$ à t , les joueurs **0, 2, 4, 8, 10** font A à t
 $\parallel e(C) \leftarrow \nabla$ pendant σ ; les joueurs **2, 8** tirent 0, les autres tirent 1

2, n_{10} : quand $e(C) \leftarrow 0$, les joueurs **2, 8** font A $\parallel e(C) \leftarrow \nabla$ pendant σ ; le joueur **8** tire 0, le joueur **2** tire 1

3, n_{11} : quand $e(C) \leftarrow 0$, les joueurs **0, 4, 10** font A $\parallel e(C) \leftarrow \nabla$ pendant σ ; ils tirent 1

4, n_{100} : quand $e(C) \leftarrow 0$, le joueur **8** fait A $\parallel e(C) \leftarrow 1$, **8** acquiert C, puis fait L, $e(C) \leftarrow 0$

5, n_{101} : quand $e(C) \leftarrow 0$, le joueur **2** fait A $\parallel e(C) \leftarrow 1$, **2** acquiert C, puis fait L, $e(C) \leftarrow 0$

6, n_{110} : $e(C)$ reste à 0

7, n_{111} : σ après transition $e(C) 1 \Rightarrow 0$, les joueurs **0, 4, 10** font A $\parallel e(C) \leftarrow \nabla$ pendant σ ; ils tirent 0

.../...

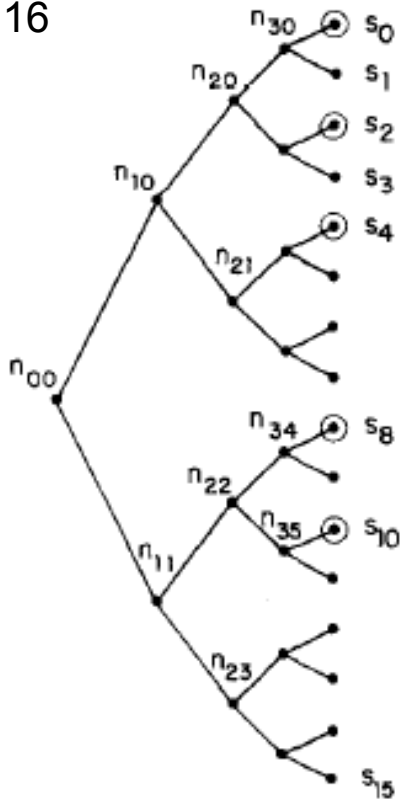
Exemple du tri binaire – Algorithme déterministe

A chaque réponse « perdu », ensemble des noms dichotomisé (itératif) :

$j \in$ moitié basse \rightarrow tenter A à nouveau à $t =$ transition $e(C) \Leftarrow 0$

$j \in$ moitié haute \rightarrow tenter A à nouveau à $t + \sigma$ et attendre transition $e(C) \Leftarrow 0$

$n = 16$



1, n_{00} : $e(C) = 0$ à t , les joueurs **0, 2, 4, 8, 10** font A à $t \parallel e(C) \Leftarrow \nabla$ pendant σ

2, n_{10} : quand $e(C) \Leftarrow 0$, les joueurs **0, 2, 4** font A $\parallel e(C) \Leftarrow \nabla$ pendant σ

3, n_{11} : quand $e(C) \Leftarrow 0$, les joueurs **8, 10** font A $\parallel e(C) \Leftarrow \nabla$ pendant σ

4, n_{20} : quand $e(C) \Leftarrow 0$, les joueurs **0, 2** font A $\parallel e(C) \Leftarrow \nabla$ pendant σ

5, n_{21} : quand $e(C) \Leftarrow 0$, le joueur **4** fait A $\parallel e(C) \Leftarrow 1$, **4** acquiert C, puis fait L, $e(C) \Leftarrow 0$

6, n_{30} : quand $e(C) \Leftarrow 0$, le joueur **0** fait A $\parallel e(C) \Leftarrow 1$, **0** acquiert C, puis fait L, $e(C) \Leftarrow 0$

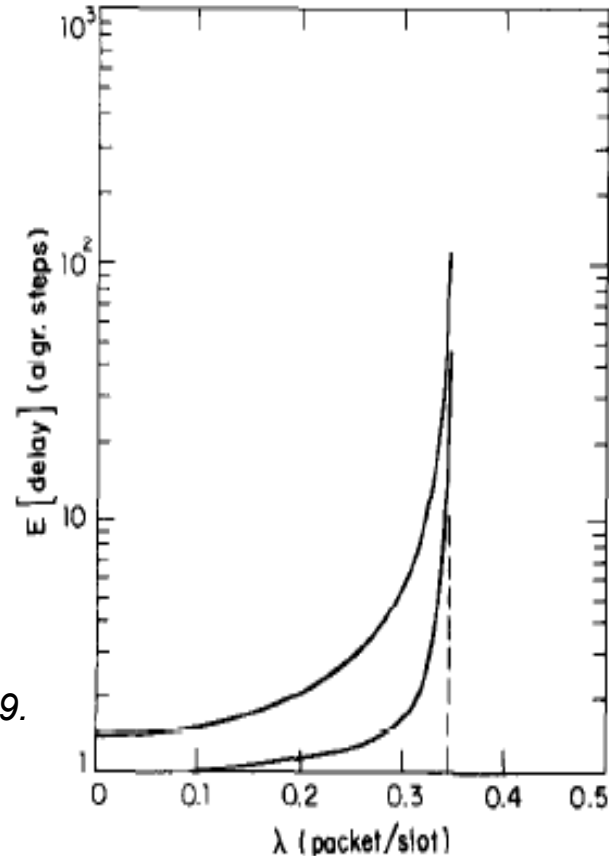
.../...

Tri binaire déterministe – Analyse markovienne des délais

δ : delay for completing a tree search (acquiring C) as a function of offered load λ (activation rate of A)

$$E\{\delta\} \leq \frac{6.05\lambda(c-2.88\lambda)}{1-(2.88\lambda)^2} + \frac{1.05(c-2.88\lambda)^2}{(1-(2.88\lambda)^2)\bar{l}_t(\lambda)} + 0.321$$

$$E\{\delta\} \geq \max \left[\frac{0.72\lambda}{1-(2.88\lambda)^2} + \frac{1}{2} \bar{l}_t(\lambda), \bar{l}_t(\lambda) \right].$$



J. I. Capetanakis, "Tree algorithms for packet broadcast channels," *IEEE Trans. Inform. Theory*, vol. IT-25, 505-515, 1979.

B. S. Tsybakov and V. A. Mikhailov, "Free synchronous packet access in a broadcast channel with feedback," *Probl. Inform. Transmission*, vol.14, pp. 259-280, 1979.

Fig. 5. Upper and lower bounds to average delay versus average arrival rate for binary tree/Poisson source system. Note: An algorithmic step equals one round-trip delay plus two slots.

Tri binaire déterministe – Analyse pires cas des délais

(analyse combinatoire)

$$\xi_k^t = \begin{cases} \frac{m \left\lceil \log_m \left(m \left\lfloor \frac{k}{2} \right\rfloor \right) \right\rceil - 1}{m-1} + m \left\lfloor \frac{k}{2} \right\rfloor \left\lceil \log_m \left(\frac{t}{m \left\lfloor \frac{k}{2} \right\rfloor} \right) \right\rceil - (k - m \left\lfloor \frac{k}{2} \right\rfloor) & \text{si } k \in \{2, \dots, t\}, \\ 0 & \text{si } k = 1, \\ 1 & \text{si } k = 0. \end{cases}$$

$$t = m^n, \quad m \in \mathbb{N}^* \setminus \{1\}, \quad n \in \mathbb{N}^*.$$

$\sigma \xi_k^t$: worst-case delay for isolating k leaves
(completing a tree search // acquiring C)
in a t-leaf balanced m-ary tree

J.-F. Hermant, G. Le Lann, "A protocol and Correctness Proofs for Real-Time High-Performance Broadcast Networks", 18th IEEE International Conference on Distributed Computing Systems (ICDCS 98), Amsterdam, The Netherlands, 26-29 May 1998, pp. 360-369.

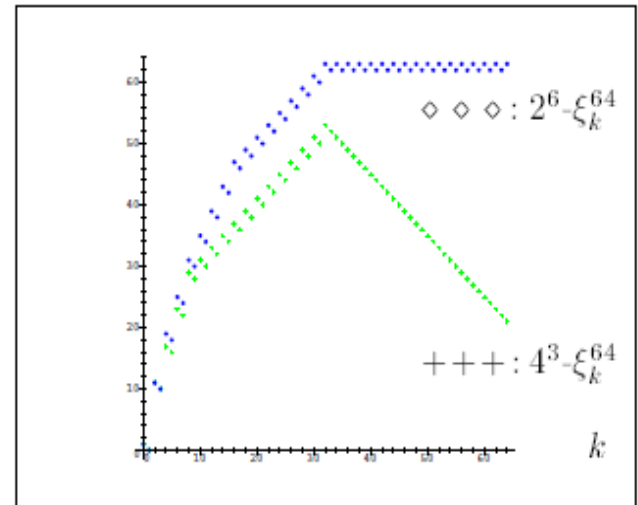


Fig. 2: Worst-case search times
for 64-leaf balanced binary and quaternary trees

Exclusion mutuelle distribuée et canaux de communication en accès multiple

La grande famille des protocoles CSMA

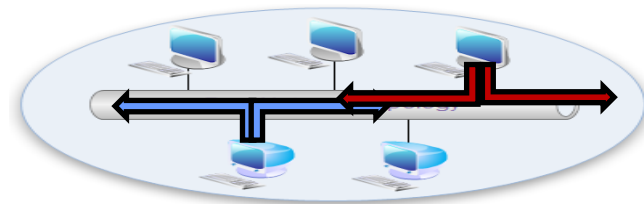
70's Palo Alto

CSMA/CD et la naissance de Ethernet (probabiliste)

Messages transmis sur câbles bidirectionnels

CSMA: Carrier Sense Multiple Access
CA: collision avoidance
CD: collision detection
CR: collision resolution

Action A : carrier sense (CS)*
si $e(C) \neq 0$ alors attendre $e(C) \leftarrow 0$
si $e(C) = 0$ alors commencer émission
et poursuivre CS
tantque $e(C) \neq \nabla$ poursuivre/terminer émission
sinon arrêter et faire BEB(d_a) *test de porteuse



$n \leq 1024$

BEB: binary exponential backoff
= tri binaire probabiliste

CD correspond à $e(C) = \nabla$. BEB (d_a) ? Tirage uniforme d'entiers sur l'intervalle $[0, 2^a - 1]$, $a = nb$ de tentatives A infructueuses \rightarrow entier $b \rightarrow d_a = b \sigma$
 \rightarrow choix aléatoire de l'une parmi 2^a feuilles d'un arbre binaire

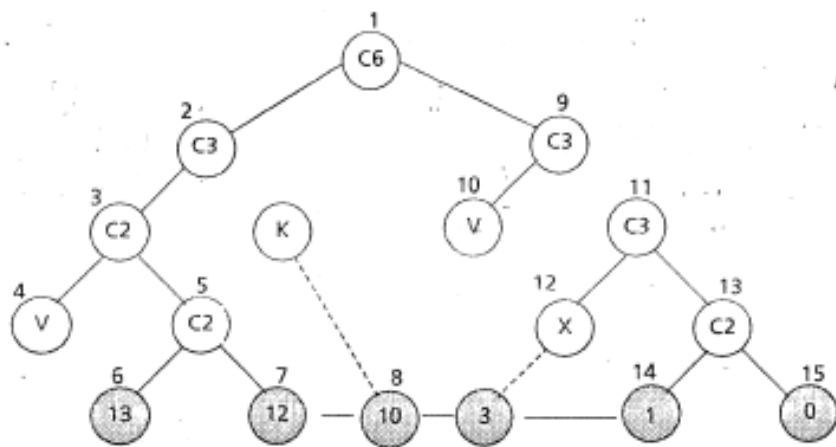
Exclusion mutuelle distribuée et canaux de communication en accès multiple

La grande famille des protocoles CSMA

80's Rocquencourt

CSMA/CD/DCR et la naissance de Ethernet déterministe

CSMA: Carrier Sense Multiple Access
 CA: collision avoidance
 CD: collision detection
 CR: collision resolution



les feuilles contiennent les valeurs d'index associés aux messages transmis

Figure 1
 Diagramme temporel de la résolution de collision (mode général)

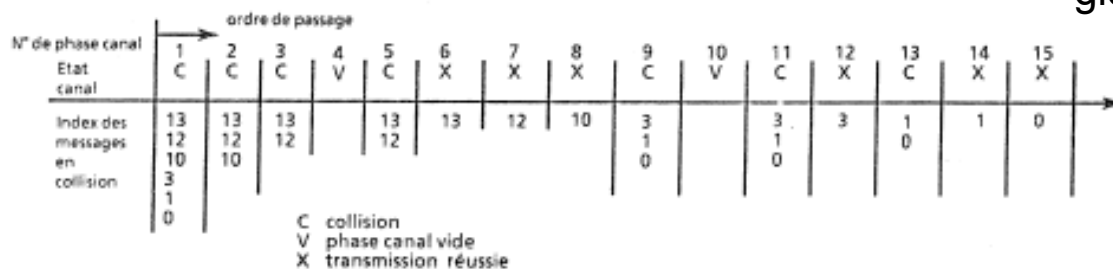
DCR: deterministic collision resolution = tri binaire déterministe

Exemple numérique

σ : ch_slot time = 40 μ s

durée utilisation C : durée message = 460 μ s

→ taux utilisation = 0,92 avec taux global d'arrivées = 2 messages/ms



Les protocoles CSMA à tris arborescents sont non bloquants

On vérifie aisément qu'un joueur k défaillant (qui s'arrête) ne peut bloquer le jeu.

- Soit k est inactif → « too bad » for k
- Soit k est actif, mais n'a pas encore acquis C → « too bad » for k
- Soit k est actif et possède C
 - sur défaillance de k , transition d'état de C : $1 \Rightarrow 0$ (arrêt de transmission)
 - les autres joueurs détectent la libération de C , et poursuivent le jeu (comme si k avait fait L)

Cette propriété n'est pas obtenue aussi aisément dans le cas de ressource C passive.

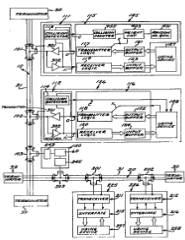
Les algorithmes à verrouillage/sémaphores un peu trop « simples » sont bloquants.

[54] MULTIPOINT DATA COMMUNICATION SYSTEM WITH COLLISION DETECTION
 [75] Inventors: Robert M. Metcalfe, Woodside; David R. Boggs, Charles P. Thacker, both of Palo Alto; Butler W. Lamson, Portola Valley, all of Calif.
 [73] Assignee: Xerox Corporation, Stamford, Conn.
 [21] Appl. No.: 563,741
 [22] Filed: Mar. 31, 1975
 [51] Int. Cl.: H04Q 9/00
 [52] U.S. Cl.: 340/147 R; 340/163; 340/346
 [58] Field of Search: 340/147 R, 147 LP, 163, 340/151, 152 R, 152 T, 346, 323/22, 57, 178/558 A; 179/2 DDP, 15 AS, 17 B, 30, 343/177, 180
 [56] References Cited
 U.S. PATENT DOCUMENTS
 3,350,687 10/1967 Gabrielson et al. 340/163
 3,496,293 2/1970 Avroy et al. 340/346
 3,509,538 4/1970 Holden et al. 340/147 R
 3,723,870 4/1973 Felchak 340/152 R
 3,742,450 6/1973 Waller 340/147 R
 3,743,460 7/1973 Trumble 352/2
 3,819,932 6/1974 Auer et al. 340/163
 3,825,897 7/1974 Lawton 340/147 R
 3,842,116 3/1976 Ferguson 352/2
 3,967,059 6/1976 Moore et al. 178/58 A

FOREIGN PATENT DOCUMENTS
 1,188,476 10/1969 United Kingdom.
 1,365,838 9/1974 United Kingdom.
 1,382,133 1/1975 United Kingdom.
 Primary Examiner—John W. Caldwell, Sr.
 Assistant Examiner—James J. Croody
 Attorney, Agent, or Firm—J. E. Deek, T. J. Anderson, B. P. Smith

ABSTRACT
 Apparatus for enabling communications between two

22 Claims, 7 Drawing Figures



or more data processing stations comprising a communication cable arranged in branched segments including taps distributed thereover. Tied to each tap is a transceiver which on the other side connects to an associated interface stage. Each transceiver includes, in addition to the usual transmitter and receiver sections, a gate which compares the data from the interface stage with the data on the cable and indicates whether such are equal. Should such be unequal, an interference between the transceiver and the cable is indicated, disabling the associated transmitter section. Each interface stage tied to such transceiver also includes an input and an output buffer on the other end thereof interfacing with a using device, such input and output buffers storing both the incoming and outgoing data. The output buffer is connected to a clock-driven shift register which converts the buffered data to a serial stream, feeds such data to a phase encoder, which then connects to the transmitter or driver section of the transceiver. The input buffer is loaded by an input shift register which derives its clock from a phase decoder, the shift register and the phase decoder both connecting to the receiver section. When the station is to start transmitting, the phase decoder detects the presence of other transmissions on the cable and delays the output shift register until no other transmissions are sensed. Once a transmission has begun, if interference is detected and the transmitter section is disabled, a random number generator is used to select an interval of time at the completion of which the next attempted transmission will take place. Concurrently, a counter counts the number of interferences, or collisions, which recur in the attempted transmissions of one data packet and weights the mean of the random number generator accordingly. The input shift register is also connected to an address decoder which enables data transfer to the input buffer only during those times when the data is preceded by an appropriate address.

[54] PROCESS AND DEVICE FOR THE TRANSMISSION OF MESSAGES BETWEEN DIFFERENT STATIONS THROUGH A LOCATION DISTRIBUTION NETWORK
 [75] Inventors: Gérard Le Lann, Paris; Pierre Roila, Les Ulis, both of France
 [73] Assignee: Inria Institute National de Recherche en Informatique et en Automatique, Le Chesnay, France

[21] Appl. No.: 820,254
 [22] Filed: Nov. 5, 1985
 [30] Foreign Application Priority Data: Nov. 7, 1984 [FR] France 84 16957
 [51] Int. Cl.: H04Q 9/00
 [52] U.S. Cl.: 370/35; 340/825.5
 [58] Field of Search: 340/825.5; 370/35, 37

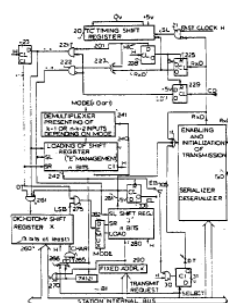
References Cited
 U.S. PATENT DOCUMENTS
 4,199,661 4/1980 White et al. 340/825.5
 4,506,361 3/1981 Kume 340/825.5
 4,512,027 4/1983 Borrillo et al. 375/97
 4,511,238 7/1983 Rawson et al. 370/85
 4,593,282 6/1984 Acampora et al. 370/85
 4,598,285 7/1984 Holton 340/825.5
 4,624,311 12/1984 Milling 370/85
 4,630,264 12/1984 Wah et al. 370/85

FOREIGN PATENT DOCUMENTS
 0088906 8/1988 European Pat. Off. 370/85
 2126848 3/1984 United Kingdom ..

OTHER PUBLICATIONS
 "The Multi-Accessing Tree Protocol", by Capetanakis IEEE Trans. on Comm. vol. Com. 27, #10, 10/10/79.
 Primary Examiner—Salvatore Cangialosi
 Attorney, Agent, or Firm—Herbert Dubno

ABSTRACT
 In a message transmission network of CSMA-CD type, each station is connected by a coupler to the common transmission channel. One or more indices is allocated to each coupler, and each coupler is provided with an automatic device capable of establishing a predetermined sequence of index sub-sets, such as a dichotomous tree. A counter of period E progresses at the rate of the end of channel phase orders. The coupler freely transmits only if its count E is at zero. When a collision appears in this state, the counter E starts from a chosen forced state, and all the couplers put their automatic devices into action, which also progress at the rate of the ends of channel phases, to establish sub-sets of said sequence and to determine, by comparison of their own indices with said sub-sets, whether they have in any given channel phase again obtained the right to transmit on said common channel. The couplers will subsequently be able to transmit one by one without collision.

29 Claims, 11 Drawing Sheets



Patented by INRIA at the request of the French Navy. Fielded in various places:

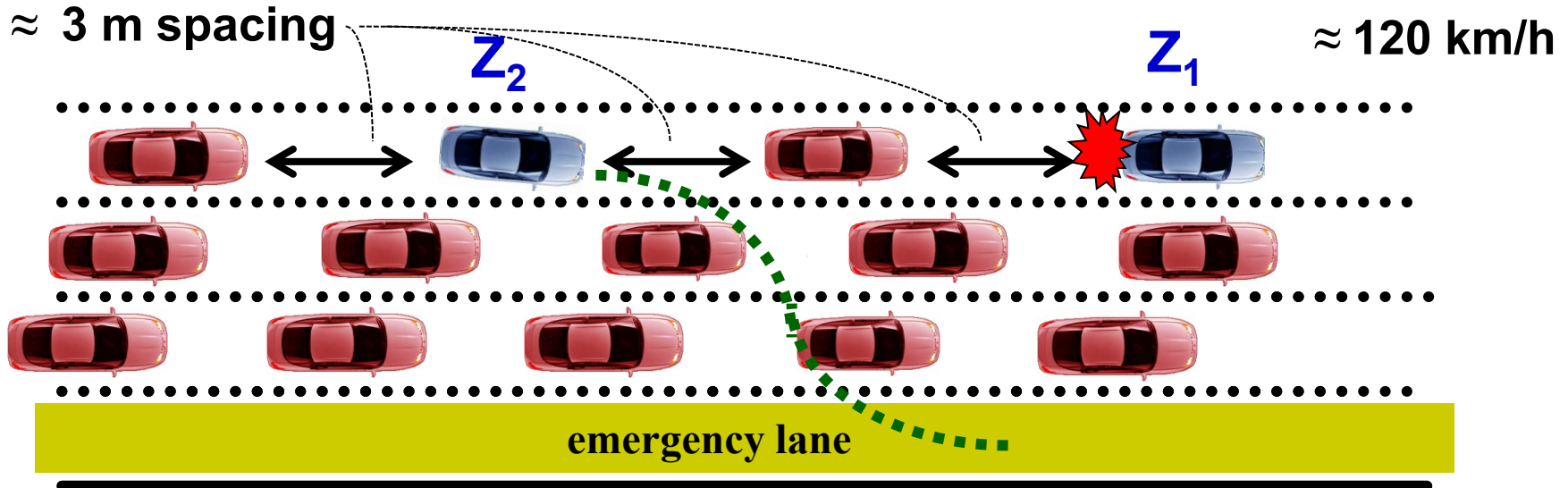
- Nuclear aircraft carrier Charles-de-Gaulle
- Frigates and submarines
- Industrial sites (Saint-Gobain, ...)
- European spatial base (Ariane launcher) in Kourou (French Guyana)

Mais noms et n sont connus → systèmes distribués « fermés »
 Domaine « fermé » : les algorithmes distribués « classiques »

Si ensemble des joueurs/processus est ouvert
 → noms ≡ inconnus, n ≡ inconnu

→ Problèmes ouverts (domaine quasi-vierge) :
 algorithmes pour systèmes cyberphysiques !

Open problems: safety-critical inter-vehicular communications (ad hoc networks of driverless/automated cars)



shared contention-prone resource C: a multiaccess radio channel

How to avoid (collective) collisions whenever risk-prone maneuvers must be undertaken. Examples:

- Z_1 decelerates abruptly
- Z_2 needs to reach the emergency lane asap

radio → câbles → radio
La grande famille des protocoles CSMA

60's Hawaiï Alohonet
 Static / radio || CSMA

70's Ethernet
 Static / câbles
 || CSMA-CD

80's Sensor Nets
 Static / radio || CSMA-CA

90's MANETs
 Mobile ad hoc nets
 Mobile / radio || CSMA-CA

00's VANETs
 Vehicular ad hoc nets
 Mobile / radio || CSMA-CA

TCP/IP et la naissance de Internet

Vint Cerf
 2012 : VP Google

Robert Metcalfe
 2012 : Polaris Venture Partners et Prof. of Innovations UTexas at Austin

Cyber Nets

Internet des Objets

CyberPhysical Nets

Défis scientifiques (retombées concrètes majeures) :
 renommage des joueurs/processus ???,
 CSMA-??? déterministe, etc.