

*Projet ADP**Algorithmes distribués et protocoles**Rennes*

THÈME 1C

 *Rapport  
d'Activité*

2000



## Table des matières

<b>1</b>	<b>Composition de l'équipe</b>	<b>3</b>
<b>2</b>	<b>Présentation et objectifs généraux</b>	<b>4</b>
<b>3</b>	<b>Fondements scientifiques</b>	<b>5</b>
3.1	Détection de propriétés . . . . .	5
3.2	Systèmes transactionnels dupliqués fondés sur la communication de groupes . . .	7
3.3	Points de contrôle et retour arrière . . . . .	8
3.4	Les problèmes d'accord dans les systèmes répartis . . . . .	11
3.5	Aide à la conception des applications multimédias . . . . .	13
<b>4</b>	<b>Domaines d'applications</b>	<b>15</b>
4.1	Panorama . . . . .	15
4.2	Le concept de communication de groupe . . . . .	15
4.3	Utilisations considérées . . . . .	16
<b>5</b>	<b>Logiciels</b>	<b>16</b>
5.1	EDEN :un service de communication de groupe . . . . .	16
<b>6</b>	<b>Résultats nouveaux</b>	<b>18</b>
6.1	Détection de propriétés . . . . .	18
6.2	Systèmes transactionnels dupliqués fondés sur la communication de groupe . . .	19
6.3	Points de contrôle et retour arrière . . . . .	20
6.4	Les problèmes d'accord dans les systèmes répartis . . . . .	21
6.5	Aide à la conception des applications multimedia . . . . .	26
<b>7</b>	<b>Contrats industriels (nationaux, européens et internationaux)</b>	<b>27</b>
7.1	Contrat FT R&D sur les contraintes non-fonctionnelles . . . . .	27
7.2	Contrat RÉUTEL 2000 . . . . .	27
<b>8</b>	<b>Actions régionales, nationales et internationales</b>	<b>28</b>
8.1	Actions nationales . . . . .	28
8.1.1	GDR ARP(Architecture, Réseaux et Parallélisme) du CNRS . . . . .	28
8.2	Actions européennes . . . . .	28
8.2.1	Cabernet . . . . .	28
8.2.2	Italie . . . . .	28
8.3	Actions internationales . . . . .	28
8.3.1	Amériques . . . . .	28
<b>9</b>	<b>Diffusion de résultats</b>	<b>29</b>
9.1	Actions d'enseignement . . . . .	29
9.2	Présentation de travaux . . . . .	30
9.3	Animation de la communauté scientifique . . . . .	31

9.4	Coopération et échanges internationaux . . . . .	32
<b>10</b>	<b>Bibliographie</b>	<b>32</b>

---

*ADP est un projet commun à l'INRIA, au CNRS et à l'université de Rennes I.*

## 1 Composition de l'équipe

### Responsable scientifique

Michel Raynal [professeur, Université de Rennes 1]

### Assistante de projet

Maryse Auffray [Adjoint Administratif, INRIA]

### Personnel INRIA

Michel Hurfin [CR]

### Personnel CNRS

Emmanuelle Anceaume [CR]

### Personnel Université de Rennes I

Jean-Michel Hélary [professeur]

Philippe Ingels [maître de conférences]

Achour Mostéfaoui [maître de conférences]

### Collaborateur extérieur

Jean-Pierre Le Narzul [maître de conférences à l'ENSTBr]

### Chercheurs invités

Giovanna Melideo [doctorant, Université de Rome, Italie, de décembre 1999 à février 2000]

Roberto Baldoni [professeur, Université de Rome, Italie, de décembre 1999 à février 2000]

Vijaykumar Krishnaswamy [professeur, Georgia Tech, College of Computing, USA, en février 2000]

Hamida Seba [maître de conférences, Institut National d'Informatique, Algérie, en mai 2000]

Manuel Romero [professeur, université du Mexique, en juin-juillet 2000]

Sergio Rajsbaum [professeur, université du Mexique, en juin 2000]

Christof Fetzer [Chercheur, AT&T Labs Research, USA, en juillet 2000]

Roman Vitenberg [doctorant, institut de technologie Technion, Israel, en septembre-octobre 2000]

### **Chercheurs doctorants**

Erwan Demairy [bourse Menesr]

Udo Fritzke [bourse du gouvernement brésilien]

Fabiola Greve [bourse du gouvernement brésilien]

Eric Mourgaya [bourse INRIA]

Matthieu Roy [allocation couplée]

Frédéric Tronel [allocation couplée]

### **Chercheurs post-doctorant**

Yun Wang [post doctorant Réutel d'octobre 1999 à octobre 2000]

Francisco Brasileiro [post doctorant Réutel d'octobre 1999 à octobre 2000]

## **2 Présentation et objectifs généraux**

L'évolution de l'informatique a été marquée ces dernières années par le développement sans précédent des réseaux qui permettent l'interconnexion de machines (il suffit de penser à INTERNET) ainsi que par l'apparition des machines dites "massivement parallèles". Ceci continue de susciter une intense activité de recherche autour des *applications* et des *systèmes répartis* et de leur fondement que constitue l'*algorithmique répartie*.

Si tous ces efforts commencent aujourd'hui à porter leurs fruits (des systèmes répartis existent), la somme des problèmes ouverts n'en demeure pas moins importante. Ainsi, beaucoup de travail, tant du point de vue théorique que pratique, reste à faire dans des domaines tels que : la sécurité, la cohérence des données réparties, les environnements de programmation répartis, la tolérance aux défaillances, etc.

Le projet ADP poursuit, dans ce contexte, deux objectifs complémentaires. Le premier réside dans la compréhension des structures "profondes" des calculs et des systèmes répartis.

Il s'agit là d'une recherche à caractère fondamental, explicitement axée (1) sur l'observation et l'analyse des exécutions réparties (détection de propriétés, définition de points de contrôle, etc.) (2) sur la synchronisation d'activités et la gestion cohérente de données réparties et (3) sur la prise en compte de contraintes non fonctionnelles telles que la tolérance aux défaillances et les contraintes temps-réel.

Le deuxième objectif, qui trouve ses fondements dans les résultats précédents et qui en aiguille la problématique, consiste en l'étude de problèmes concrets, à savoir la réalisation d'un noyau réparti de communication fiable et la conception de protocoles pour les applications multimedia.

Comme indiqué précédemment, ces deux objectifs (théorie et pratique) sont étroitement imbriqués dans la philosophie du projet ADP et, en conséquence, y sont vus comme les deux faces indissociables d'une même entité, à savoir l'algorithmique du contexte réparti.

### 3 Fondements scientifiques

La compréhension des fondements des algorithmes répartis passe par la recherche des concepts sous-jacents et de nouveaux paradigmes. Elle doit permettre de concevoir de nouveaux algorithmes répartis, fournissant des services de base tels que le contrôle de la compétition pour les ressources, la détection du passage d'un calcul dans un ensemble d'états particuliers, la collecte d'informations réparties, l'incidence de la causalité, la tolérance aux défaillances, etc. Les travaux actuels s'articulent autour des thèmes ci-dessous.

#### 3.1 Détection de propriétés

**Mots clés :** algorithme réparti, causalité, détection de propriétés, détection décentralisée.

**Résumé :** *La détection de propriétés dans les exécutions réparties est un problème fondamental qui se pose pour le concepteur d'environnements de programmation répartie ainsi que pour le contrôle des programmes répartis. Malheureusement, du fait de la structure de l'ensemble des états globaux cohérents d'une exécution répartie, la détection de propriétés est un problème NP-complet.*

La détection de propriétés dans les exécutions réparties est un problème fondamental qui se pose pour le concepteur d'environnements de programmation répartie ainsi que pour le contrôle des programmes répartis. Pour le concepteur d'environnements, il s'agit essentiellement de pouvoir fournir un outil de compréhension de ces programmes. Par exemple, lors d'une séance de débogage, il peut demander si une somme de variables réparties représentant une ressource bornée est supérieure à une constante fixée, s'il existe une conjonction d'états locaux qui forment un état global cohérent tel que chacun de ces états locaux respecte une propriété donnée. Pour le contrôle de programmes répartis, la détection de comportement est un aspect du problème. Ainsi avant de mettre en oeuvre un algorithme de résolution d'interblocage, il faut savoir détecter si un sous-ensemble de processus est en situation d'interblocage.

Dans son cas le plus général, la détection de propriétés est un problème NP-complet. L'observation du comportement d'une exécution par des observateurs répartis pouvant donner

lieu à des observations différentes pour un même comportement, il est nécessaire de générer toutes les observations possibles [SM94]. Ceci est dû à la structure de l'ensemble des états globaux cohérents d'une exécution répartie qui au pire, a une taille exponentielle par rapport au nombre de processus mis en jeu. Pour ne pas être limité par ces caractéristiques dans la conception d'algorithmes de détection, il est alors essentiel d'étudier ce problème de manière moins générale, soit en restreignant les types de propriétés que l'on souhaite détecter, soit en restreignant la structure sur laquelle sont détectées les propriétés.

Historiquement, c'est la restriction aux types de propriétés qui a d'abord été employée à travers la distinction *propriétés stables/propriétés instables*. Une propriété est stable si une fois vérifiée dans une exécution répartie, elle le demeure. On peut montrer qu'il existe des algorithmes de coût polynomial pour résoudre le problème de leur détection. Toutefois le problème demeure car la classe des propriétés stables ne rassemble qu'un nombre limité de cas particuliers.

Toujours dans le domaine de la restriction au type de propriétés, une voie plus prometteuse a été récemment étudiée. Il s'agit de la limitation syntaxique dans l'expression même des propriétés. C'est dans cette approche que s'inscrivent les travaux de V. Garg concernant les conjonctions de prédicats locaux [GW94] et les prédicats linéaires ou semi-linéaires [CG95], ainsi que les travaux de Stoller et Schneider sur la décomposition d'expressions en conjonctions.

Au sein de l'équipe ADP, nous nous intéressons plus particulièrement à la restriction de la structure sur laquelle sont détectées les propriétés. Nous avons présenté le problème de la détection de propriétés comme une variante d'un problème de reconnaissance de langage. Le calcul réparti est modélisé sous la forme d'un graphe orienté acyclique. Les sommets du graphe correspondent à des états (locaux ou globaux). Les arcs du graphe sont le reflet des relations de précedence causale qui existent entre les états. Des prédicats de bases (locaux ou globaux) sont exprimés sous forme d'expressions booléennes portant sur les valeurs des variables gérées par les processus de l'application. Ces prédicats de base sont identifiés de manière unique par une lettre (ou symbole). Par définition, un symbole est associé à un sommet du graphe si l'état correspondant satisfait le prédicat de base ainsi identifié. L'ensemble des symboles spécifiés constitue un alphabet. Nous appelons *motif* tout langage défini sur cet alphabet. Le rôle d'un algorithme de détection est de comparer un ensemble de motifs spécifiés avec l'ensemble des motifs observés au cours de l'exécution. Cette comparaison se fait en visitant les sommets du graphe selon la stratégie du tri topologique. Une règle de satisfaction (ou opérateur modal) précise les modalités de comparaison. En choisissant le type de graphe (états globaux, états locaux ou événements) nous pouvons régler le niveau de complexité de l'algorithme. Cette stratégie nous a notamment permis de reformuler différents algorithmes de détection de propriétés dans un cadre général unique ainsi que d'en proposer de nouveaux [BFR96].

- 
- [SM94] R. SCHWARZ, F. MATTERN, «Detecting Causal Relationships in Distributed Computations : In Search of the Holy Grail», *Distributed Computing* 7, 3, 1994, p. 149–174.
- [GW94] V. K. GARG, B. WALDECKER, «Detection of weak unstable predicates in distributed programs», *IEEE Trans. on Parallel and Distributed Systems* 5, 3, 1994, p. 299–307.
- [CG95] C. M. CHASE, V. K. GARG, «Efficient Detection of Restricted Classes of Global Predicates», in : *9th International Workshop on Distributed Algorithms (WDAG'95)*, J.-M. Hélary, M. Raynal (éditeurs), Springer-Verlag, LNCS 972, p. 303–317, 1995.
- [BFR96] O. BABAOĞLU, E. FROMENTIN, M. RAYNAL, «A Unified Framework for the Specification and Run-

## 3.2 Systèmes transactionnels dupliqués fondés sur la communication de groupes

**Mots clés :** transaction, duplication, diffusion, tolérance aux défaillances.

**Résumé :** *La duplication augmente la disponibilité des données et permet d'exécuter des transactions même en cas de défaillances. Dans les systèmes transactionnels dupliqués, au delà d'assurer l'isolation des transactions, la duplication doit être rendue transparente à l'utilisateur, et cela tout en tolérant les défaillances. En raison de la distribution des données, on doit également mettre en œuvre une validation atomique répartie pour les transactions. Par ailleurs, contrairement aux systèmes non-dupliqués, les transactions impliquent des ressources qui sont multipliées sur un système réparti et par conséquent, elles sont plus sensibles aux problèmes d'interblocage. Le traitement de ces problèmes s'avère d'autant plus difficile que l'on se place dans le cadre des systèmes répartis asynchrones.*

*Les primitives de diffusion peuvent y jouer un rôle déterminant dans la mesure où elles rendent les défaillances et l'asynchronisme transparents à l'application. Enfin, elles procurent des propriétés d'ordre qui rendent l'implémentation du contrôle de la duplication plus simple.*

La duplication de données a été utilisée depuis longtemps pour assurer que les données accédées par des transactions demeurent disponibles en dépit des défaillances. De nombreux protocoles de contrôle de duplication qui assurent la "sérialisabilité à une copie" (one-copy serializability) des transactions ont été proposés. Ces protocoles tolèrent différentes situations de défaillances (défaillances de processus ou partitionnement du réseau, par exemple), mais d'autre part, ils négligent d'autres aspects liés à la vivacité des transactions. En l'occurrence, ni l'interblocage de transactions ni la validation atomique non-bloquante de celles-ci ne sont traités par ces protocoles. Certes, cela se doit en grande partie au fait que l'on ne pose que des hypothèses faibles concernant les propriétés des systèmes de communication sous-jacents. Souvent, on ne considère que des primitives de communication point-à-point, voire, sans garantie de fiabilité de la livraison des messages. La duplication s'avère encore plus difficile à gérer lorsque l'on se place dans le cadre des systèmes répartis asynchrones, dans lesquels aucune hypothèse sur les vitesses des communications ou des processeurs n'est faite. Il est connu que des problèmes d'accord concernant la validation atomique et le verrouillage réparti sont impossibles à résoudre dans les systèmes asynchrones.

À l'opposé des primitives de communication considérées par les protocoles traditionnels, les systèmes de communication de groupes procurent, entre autres, des primitives de diffusion de messages avec des garanties fortes de fiabilité et d'ordre de livraison. Par ailleurs, il existe des protocoles qui implémentent ces primitives de diffusion dans les systèmes asynchrones. Avec l'utilisation de tels services de diffusion au sein des systèmes transactionnels dupliqués, les aspects concernant le traitement d'interblocage et la validation atomique ont pu être facilement incorporés aux protocoles de contrôle de duplication, et cela même dans un cadre asynchrone.

Des travaux récents ont proposé cette intégration de la communication de groupe avec la gestion de la duplication. Pourtant, on y considère un modèle de duplication *complète*, où chaque processus du système réparti possède une copie complète d'une seule base de données. Contrairement à cette duplication complète, nous considérons un modèle de duplication *partielle* où chaque donnée peut avoir des copies sur un nombre quelconque de processus. Plus précisément, nous considérons que le système réparti est structuré en plusieurs groupes de processus, chaque groupe réalisant la duplication d'un objet particulier. Les transactions peuvent donc accéder à des objets sur différents groupes de processus. En conséquence, pour assurer les propriétés transactionnelles, on doit mettre en œuvre un système de gestion de transactions réparti.

Dans la duplication complète on utilise un cas particulier de diffusion, dans laquelle un message est destiné à tous les processus du système réparti ("broadcast"). En revanche, pour gérer de manière efficace la duplication partielle, les messages doivent être diffusés vers des groupes de processus, différents messages pouvant avoir des ensembles de groupes destinataires distincts ("multicast"). Les primitives de *Diffusion fiable uniforme* et de *Diffusion fiable totalement ordonnée* (ou *Diffusion atomique*) vers plusieurs groupes ont ces caractéristiques. La première primitive implémente une propriété de fiabilité qui garantit qu'un message est inéluctablement délivré par les processus corrects des groupes destinataires si l'émetteur est correct, ou si au moins l'un des processus a délivré ce message. La deuxième met en œuvre cette fiabilité et en plus assure que les messages sont délivrés en respectant un ordre total.

Pour faire face aux systèmes répartis en grande échelle, il est important par ailleurs que les protocoles de diffusion réalisent la propriété de *Minimalité*. Cette propriété garantit que seuls le processus émetteur et les processus des groupes destinataires d'un message sont impliqués dans sa diffusion et sa livraison.

### 3.3 Points de contrôle et retour arrière

**Mots clés :** algorithme réparti, causalité, points de contrôle, cohérence, tolérance aux défaillances.

**Résumé :** *Déterminer des points de contrôle globaux cohérents est une tâche importante qui trouve des applications tant dans le domaine de la tolérance aux défaillances (points de reprise à partir desquels un calcul peut être relancé après défaillance) que dans la détection de propriétés des exécutions réparties (coupes cohérentes). De telles déterminations sont non triviales dans le cadre des systèmes répartis asynchrones. La théorie des Z-chemins de Netzer et Xu énonce qu'un tel chemin entre deux points de contrôle locaux révèle une dépendance qui leur interdit d'appartenir à un même point de contrôle global cohérent.*

Lors de l'exécution répartie d'une application, un point de contrôle global est un ensemble de points de contrôle locaux (états locaux), un par processus participant à l'exécution. La détermination d'un point de contrôle global cohérent est un problème important dans de nombreux domaines concernés par les applications réparties (résistance aux défaillances, mise au point

répartie, détection de propriétés, etc.). De nombreux protocoles ont été proposés pour déterminer des points de contrôle locaux afin qu'ils forment des points de contrôle globaux cohérents [BHMR95,EJW96]. Si les points de contrôle locaux sont sélectionnés de manière non coordonnée, il est possible que ceux-ci ne puissent former aucun point de contrôle global cohérent (ce risque est connu sous le nom d'*effet domino* [Ran75]). Une forme de coordination est donc nécessaire si l'on veut éviter - ou diminuer - cet effet. Les techniques de calcul d'état global cohérent introduites par Chandy-Lamport [CL85] sont basées sur une coordination explicite utilisant des messages de contrôle supplémentaires. En général, cette forme de coordination impose à tous les processus de prendre des points de contrôle locaux lorsque l'un d'entre eux décide d'en prendre un, même en l'absence de communication due à l'exécution répartie. Une autre forme de coordination, implicite, utilise des informations de contrôle véhiculées par les messages de l'application (technique de *piggybacking*). Pour cette raison, elle est connue sous le nom de *coordination induite par les communications* (communication-induced checkpointing). Dans cette approche, des points de contrôle locaux sont sélectionnés de manière non coordonnée (points de contrôle *de base*) et le protocole impose des points de contrôle locaux supplémentaires, appelés *points de contrôle forcés*, de manière à assurer la progression des points de contrôle globaux cohérents. Les points de contrôle forcés sont pris sur la base des informations de contrôle véhiculées par les messages de l'application. C'est cette approche qui a été principalement étudiée dans notre équipe depuis l'année 1995, bien que nous ayons abordé récemment l'approche non coordonnée.

**Études dans le cadre de la notion classique cohérence** Le fait que deux points de contrôle locaux ne soient pas liés causalement constitue une condition nécessaire pour appartenir à un même point de contrôle global cohérent. Malheureusement, cette condition n'est pas suffisante. Les points de contrôle locaux peuvent avoir des dépendances cachées (c'est-à-dire non captables par un mécanisme d'estampillage) qui les empêchent de participer au même point de contrôle global cohérent. Afin de capter l'ensemble des dépendances liant les points de contrôle locaux, Netzer et Xu [NX95] ont introduit la notion de *Z-chemin* entre points de contrôle locaux. Ils ont de plus démontré le théorème fondamental suivant : *un ensemble quelconque de points de contrôle locaux peut être étendu pour former un point de contrôle global cohérent si, et seulement si, il n'y a pas de Z-chemin connectant deux points de contrôle de cet ensemble*. Si, dans une exécution répartie, on ne considère que les points de contrôle locaux et leurs relations de dépendance, on obtient une abstraction de cette exécution répartie (cette

- 
- [BHMR95] J. BRZEZINSKI, J. HÉLARY, A. MOSTEFAOUI, M. RAYNAL, « Semantics of recovery lines for backward recovery in distributed systems », *Annales des Télécommunications* 50, 11-12, 1995, p. 874-887.
- [EJW96] E. ELNOZAHY, D. JOHNSON, Y. WANG, « A Survey of Rollback-Recovery Protocols in Message-Passing Systems », *Technical Report n° CMU-CS-96-181*, Carnegie-Mellon University, 1996.
- [Ran75] B. RANDELL, « System Structure for Software Fault-Tolerance », *IEEE Transactions on Software Engineering SE1*, 2, 1975, p. 220-232.
- [CL85] K. CHANDY, L. LAMPORT, « Distributed Snapshots: Determining Global States of Distributed Systems », *ACM TOCS*, February 1985, p. 63-75.
- [NX95] R. NETZER, J. XU, « Necessary and Sufficient Conditions for Consistent Global Snapshots », *IEEE Transactions on Parallel and Distributed Systems* 6, 2, 1995, p. 165-169.

abstraction ignore tous les états locaux qui ne sont pas des points de contrôle locaux). Une question importante est alors : “cette abstraction est-elle cohérente ?” Cette question peut être abordée dans deux contextes différents, selon la notion de cohérence particulière considérée pour l’abstraction définie par les points de contrôle.

Dans le premier cas, l’abstraction est cohérente si tout point de contrôle local appartient à au moins un point de contrôle global cohérent. En d’autres termes, il n’y a pas de point de contrôle local qui s’avère inutile du point de vue de la construction des points de contrôle globaux cohérents (du point de vue opérationnel ceci élimine l’effet domino lors de la construction des points de contrôle globaux). Netzer et Xu ont montré, au niveau de l’abstraction, qu’il n’y a pas de point de contrôle inutile si, et seulement si, aucun Z-chemin n’est un cycle (Z-cycle). Il s’agit alors de traduire cette propriété au niveau opérationnel et d’en déduire des protocoles efficaces, notamment dans le cadre de l’approche *coordination induite par les communications* (conception et analyse de protocoles de détermination de points de contrôle forcés prévenant l’occurrence de points de contrôle inutiles).

Une seconde notion de cohérence pour l’abstraction que constitue un ensemble de points de contrôle locaux réside dans l’absence de relations de dépendance cachées entre ces points de contrôle locaux. Cette notion de cohérence, plus forte que la précédente, a été proposée par Wang [Wan97]. Appelée propriété RDT (pour *Rollback-Dependency Trackability*), elle se révèle particulièrement utile pour résoudre des problèmes tels que le calcul au vol du point de contrôle global *minimal* cohérent auquel appartient un point de contrôle local donné (Dans [Wan97] d’autres exemples de problèmes dont la solution est facilitée lorsque les points de contrôle définis lors d’une exécution satisfont la propriété RDT. A titre d’exemples, les points de contrôle globaux cohérents minimaux facilitent la recherche d’erreurs logicielles et la relance des exécutions après la détection d’un interblocage). Dans le cadre de l’approche *coordination induite par les communications*, nous nous sommes intéressés à la recherche de conditions optimales sous lesquelles un processus sera obligé de prendre un point de contrôle forcé. Par là nous entendons la recherche de conditions telles que le nombre de points de contrôle forcés serait le plus faible possible. Dans une approche “théorique”, il s’agit de définir un ensemble de Z-chemins dont la suppression (par prise de points de contrôles locaux) suffit à assurer la propriété RDT. Trivialement, cette propriété est assurée si *tous* les Z-chemins sont ainsi supprimés. Mais cette solution peut s’avérer coûteuse en termes du nombre de points de contrôle forcés qu’elle induit. C’est pourquoi la recherche d’ensembles *minimaux* de Z-chemins est un problème de toute première importance. Sur le plan “pratique”, il s’agit de concevoir des protocoles “génériques” qui offrent à l’utilisateur des compromis entre taille des informations de contrôle véhiculées par les messages et nombre de points de contrôle forcés.

**Extensions de la notion de cohérence** Un examen approfondi des notions de cohérence d’états globaux (ou de points de contrôle globaux) fait apparaître que la notion usuelle, basée sur la causalité, ne représente qu’un aspect du problème. C’est pourquoi d’autres modèles de cohérence doivent être envisagés. Deux modèles, en particulier, présentent un grand intérêt pratique. Ce sont, respectivement, *l’absence de messages en transit (transitlessness)*, duale de

---

[Wan97] Y. WANG, « Consistent Global Checkpoints That Contain a Given Set of Local Checkpoints », *IEEE Transactions on Computers* 46, 4, 1997, p. 456–468.

la notion classique, et la *cohérence forte*, réunion de la cohérence classique et de l'absence de messages en transit. La cohérence classique considère des points de contrôle globaux dans lesquels tout message est enregistré comme *reçu* seulement s'il y est enregistré comme *émis*. La cohérence *transitless* considère des points de contrôle globaux dans lesquels tout message doit être enregistré comme *reçu* s'il y est enregistré comme *émis*. La cohérence forte considère des points de contrôle globaux dans lesquels tout message est enregistré comme *reçu* si, et seulement si, il y est enregistré comme *émis*. Le problème consiste alors à définir des formalismes permettant d'étendre à ces modèles de cohérence la condition nécessaire et suffisante pour qu'un ensemble de points de contrôle locaux puisse faire partie d'un point de contrôle global cohérent. Cette étude théorique sert de fondement à la conception de protocoles assurant la cohérence forte en combinant les deux techniques de points de contrôle forcés et de sauvegarde des messages (recording), ainsi qu'à l'étude des compromis entre ces deux techniques.

### 3.4 Les problèmes d'accord dans les systèmes répartis

**Mots clés :** algorithme réparti, tolérance aux défaillances, temps réel, communication de groupe, consensus.

**Résumé :** *Dans un système réparti asynchrone, il est important de pouvoir concevoir des applications tolérant les défaillances tout en garantissant le respect de contraintes temps-réel. Le concept de groupe s'avère dans ce cas particulièrement intéressant. Développer les services offerts dans un groupe en utilisant comme brique de base une solution au problème du consensus est une approche novatrice présentant de nombreux avantages. En particulier, le résultat d'impossibilité de Fischer-Lynch-Paterson peut être circonscrit à ce niveau.*

Considérer qu'un système réparti est asynchrone est une hypothèse attrayante (car plus générale et réaliste). Cependant, elle semble difficilement conciliable avec la prise en compte de critères de qualité de service tels que la tolérance aux défaillances et les contraintes temps-réel. Un des axes de recherche du projet ADP est consacré à la conception et au développement de services permettant de répondre efficacement à ces nouvelles exigences dans un environnement distribué asynchrone.

La sûreté de fonctionnement ne peut être garantie que pour des types de défaillances préalablement identifiés. Dans le cadre de cette activité, nous considérons les défaillances de type panne franche (crash) : chaque processus peut, soit s'exécuter correctement, soit s'interrompre brutalement (et définitivement) suite à une panne (ou une agression extérieure). A condition de pouvoir coordonner efficacement l'activité des processus dupliqués, le concept de groupe se révèle être une excellente technologie *middleware* pour concevoir des mécanismes de tolérance aux défaillances. Pour mettre en œuvre ce concept, il convient d'apporter des solutions efficaces à divers problèmes d'accord entre processus. En particulier, les processus appartenant à un même groupe doivent être unanimes en ce qui concerne l'identité des membres qui le composent (membership) et l'ordre dans lequel les messages qui leurs sont adressés seront délivrés (diffusion ordonnée).

Des travaux récents ont mis en évidence le lien qui existe entre les problèmes nécessitant l'obtention d'un accord (élection, diffusion ordonnée, validation atomique non-bloquante, gestion des membres d'un groupe, etc.) et le problème abstrait du consensus. De fait, ce problème élémentaire est l'objet de nombreux travaux de recherche depuis quelques années. Le résultat le plus important est malheureusement négatif [FLP85] : ce problème fondamental ne peut pas être résolu dans des systèmes asynchrones dès lors que les processus sont susceptibles de stopper prématurément leurs exécutions. Afin de surmonter ce résultat d'impossibilité, Chandra et Toueg [CT96] ont augmenté le modèle asynchrone en introduisant la notion de détecteur de défaillances. Un détecteur est associé à chaque processus et est chargé de détecter les défaillances externes. La mise en œuvre du mécanisme de détection se fait en définissant des délais de garde : un processus est suspecté s'il ne s'est pas manifesté avant l'expiration du délai de garde. Cela signifie que (1) la détection d'une défaillance réelle est généralement différée et que (2) un détecteur de défaillances peut commettre des erreurs en suspectant à tort un processus d'avoir stoppé son exécution. Chandra et Toueg définissent huit classes de détecteurs de défaillances en caractérisant chacune d'entre elles par une propriété de complétude et une propriété d'exactitude. Une propriété de complétude définit des contraintes concernant la détection des processus réellement arrêtés tandis que la propriété d'exactitude vise à limiter les suspicions erronées que peut commettre un détecteur de défaillances.

Parmi ces classes, la classe dénotée  $\diamond S$  est particulièrement intéressante puisqu'il s'agit de la classe la plus faible permettant de résoudre le consensus [CHT96]. Cette classe est caractérisée par une propriété de complétude forte (tout processus défaillant finit par être suspecté de façon permanente par tout processus correct) et une propriété d'exactitude faible inéluctable (il existe un instant à partir duquel un processus correct ne sera plus jamais suspecté par aucun processus correct). En s'appuyant sur des détecteurs de défaillances de cette classe et à condition que les canaux de communication soient fiables et qu'une majorité de processus ne subit pas de défaillances, des algorithmes déterministes permettent de résoudre le problème du consensus. Tous s'appuient sur le paradigme du coordinateur tournant mais diffèrent par leur complexité en temps et en nombre de messages. L'utilisation de tels algorithmes comme briques de base dans la résolution de problèmes d'accord, et par conséquent leur utilisation dans la gestion des groupes de processus dupliqués, est actuellement un axe de recherche important.

Comme indiqué précédemment, nous souhaitons intégrer la notion de contrainte temps-réel dans les solutions que nous proposons aux problèmes d'accord ainsi que dans les services de duplication mis en œuvre pour garantir la sûreté de fonctionnement. Prendre en compte des contraintes temps-réel nécessite "à un niveau ou à un autre" de considérer le temps physique. C'est pour cela que nous comptons remplacer le modèle "totalement" asynchrone par le modèle "asynchrone temporisé". Ce modèle réaliste est défini par la caractéristique suivante. Tout processus peut définir des délais maximaux sur les temps de transfert et de traitement. Mais comme le support est asynchrone, ces délais peuvent être violés : dans ce cas, l'application en

- 
- [FLP85] M. FISCHER, N. LYNCH, M. PATERSON, «Impossibility of Distributed Consensus with One Faulty Process», *Journal of the ACM* 32, 2, April 1985, p. 374–382.
- [CT96] T. CHANDRA, S. TOUEG, «Unreliable Failure Detectors for Reliable Distributed Systems», *Journal of the ACM* 34, 1, March 1996, p. 225–267.
- [CHT96] T. CHANDRA, V. HADZILACOS, S. TOUEG, «The Weakest Failure Detector for Solving Consensus», *Journal of the ACM* 43, 4, July 1996, p. 685–722.

est avertie (c'est la notion de *fail-awareness* décrite dans [FC96]) et réagit par un traitement d'exception approprié.

La définition de "bons" délais, dépend du support et de sa charge en "régime de croisière"; les manquements à ces délais n'apparaissent qu'en périodes d'instabilité dues à des surcharges occasionnelles ou à l'occurrence de situations imprévisibles.

### 3.5 Aide à la conception des applications multimédias

**Mots clés :** application multimedia, architecture qualité de services, temps réel, allocation de ressources.

**Résumé :** *Un nombre croissant d'applications reposent sur l'utilisation du multimedia : au sein d'une même présentation peuvent être utilisées à la fois des images animées, du son, du texte, etc. Ces applications se caractérisent par leur caractère distribué, par leur contraintes temporelles, et par la grande quantité de ressources requises. Pour tenter d'assouplir cette dernière caractéristique, un nombre croissant d'applications multimedia autorisent différents niveaux de qualité pour un même service, permettant ainsi d'adapter la qualité en fonction des ressources disponibles.*

La généralisation des applications multimedia conduit inéluctablement à modifier l'infrastructure qui répartit les ressources entre les applications. Jusqu'à l'émergence du multimedia, les exigences de qualité de service des applications effectuant des transferts d'information se limitaient le plus souvent à vérifier l'intégrité des données échangées, afin d'éviter, par exemple, la corruption des fichiers transférés. Dans ce contexte, les ressources ne sont en général pas réservées, ce qui rend imprévisible la durée d'un traitement. Ce modèle, parfois qualifié de « paresseux (*best effort*) », est acceptable pour des applications où le traitement effectivement effectué est prépondérant sur les délais. *A contrario*, le multimedia se caractérise, outre l'exigence - devenue « seulement » statistique - d'intégrité des données, par le respect de garanties temporelles, ce qui nécessite une certaine disponibilité du processeur, de la bande passante, de la mémoire, etc. [BS97].

Le premier problème à résoudre est de vérifier la *cohérence des formats des données* échangées, afin de savoir si les différents équipements utilisés sont compatibles. Comme les flux multimedia sont maintenant bien standardisés, il n'est pas nécessaire de connaître le détail des protocoles utilisés. On peut considérer que ces flux sont typés pour effectuer la vérification de cohérence des formats des données échangées.

Le second problème est posé par le respect des *contraintes temporelles*. Il s'agit de savoir si les délais et les fréquences d'arrivée des données sont compatibles avec les attentes des utilisateurs. Les contraintes temporelles sont de trois types. Les contraintes de synchronisation intra-flux expriment la régularité du flux en précisant le délai d'arrivée entre deux trames

---

[FC96] C. FETZER, F. CRISTIAN, « Fail-Awareness in Timed Asynchronous Systems », *in: Proc. of the fifteenth ACM Symposium on Principles of Distributed Computing (PODC'96)*, Philadelphie, May 1996.

[BS97] G. BLAIR, J. STEPHANI, « Open Distributed Processing and Multimedia », 1997.

d'information. Les synchronisations inter-flux précisent les délais qui peuvent séparer l'arrivée des trames sur plusieurs flux liés. Enfin, le délai de bout-en-bout concerne le délai qui s'écoule entre le passage d'une trame en un point de l'application et son arrivée en un autre point. Afin de vérifier le respect de ces contraintes, il est nécessaire de pouvoir connaître les performances des éléments employés, et de pouvoir en déduire si leur composition respecte les contraintes temporelles qui sont exprimées sur l'application.

Le comportement temporel des éléments constituant l'application ne peut être garanti que si une certaine quantité de ressources est disponible. En effet, il est possible de trouver au sein de la plate-forme d'exécution plusieurs applications multimedia qui s'exécutent en concurrence. Dans ce cas, les ressources du calculateur et du réseau ont à être partagées par toutes les applications.

Afin de permettre aux applications multimedia, très exigeantes en ressources, de s'exécuter avec le plus grand éventail possible de disponibilités en ressources, on a donné la faculté aux applications multimedia de fonctionner selon des modes plus ou moins dégradés. Dans un premier temps on a pu se satisfaire de solutions *ad hoc*, dans lesquelles les applications multimedia intégraient en leur sein des mécanismes d'adaptation de leur comportement en fonction des ressources disponibles. Mais ces solutions ont pour inconvénient d'être difficilement réutilisables et elles ne permettent pas une gestion coordonnée des ressources lorsque plusieurs applications s'exécutent.

Des solutions middleware - appelées *architectures qualité de service* - ont été proposées afin de décharger les applications de la gestion de l'adaptation. Pour profiter de leurs services, les applications doivent se conformer à un style d'architecture logicielle. Ce style spécifie entre autre les modalités d'interactions, l'interface des applications et le type des données pouvant être échangées.

Le rôle des architectures qualité de service (que nous abrégons dorénavant par AQDS) est d'effectuer la réservation et la renégociation des ressources de « bout-en-bout »<sup>[ACH95]</sup>. Ces mécanismes sont suffisants pour des applications où le serveur est lié au client directement par un réseau, comme c'est le cas souvent pour les applications de vidéo à la demande ; mais ces AQDS ne permettent pas d'appréhender correctement les liaisons de bout-en-bout d'une application qui mettent en jeu plusieurs liens de machine à machine.

Quant aux architectures qualité de service qui ne font pas d'hypothèse sur les éléments qui sont placés entre les extrémités, la complexité de leur représentation rend douteuse l'utilité à grande échelle des algorithmes proposés.

Ces architectures qualité de service définissent une infrastructure de connexion, mais ne permettent de savoir qu'à l'exécution si une application multimedia peut s'exécuter ou non. Aussi nous sommes nous intéressés à la caractérisation de ces architectures afin de permettre de déterminer *a priori* si la combinaison d'une application multimedia et d'une AQDS est, d'une part, compatible pour le format des données échangées et, d'autre part, que cette combinaison permet bien le respect des contraintes temporelles de l'application multimedia.

---

[ACH95] C. AURRECOECHA, A. CAMPBELL, L. HAUW, « A review of Quality of Service Architectures », *ACM Multimedia Systems Journal*, nov 1995.

## 4 Domaines d'applications

### 4.1 Panorama

**Mots clés :** application répartie, résistance aux défaillances, temps réel, communication de groupe, consensus, CORBA, hétérogénéité.

**Résumé :** *Suite au développement récent des technologies réseaux, la conception de services nouveaux destinés à être exécutés dans des environnements distribués hétérogènes connaît un essor important dans de nombreux domaines industriels et plus particulièrement dans le domaine des télécommunications. L'émergence récente des applications de type travail coopératif est, à ce titre, un exemple significatif.*

*Enrichir une norme telle que CORBA en spécifiant un ensemble de services destiné à assurer simultanément le respect de contraintes de sûreté de fonctionnement et de contraintes temps-réel correspond donc à une attente de la part de nombreux industriels. C'est dans cette voie que nous développons actuellement des coopérations avec des partenaires industriels.*

### 4.2 Le concept de communication de groupe

Dans un système réparti, un service de *communication de groupe* est un outil puissant qui permet de gérer les communications multi-point en rassemblant les processus qui sont destinataires des mêmes messages au sein d'une entité d'adressage unique : *le groupe*. Au niveau applicatif, un groupe est donc un ensemble de processus ayant un comportement similaire au sens où tous les membres du groupe traitent les mêmes requêtes. Un processus peut appartenir à un ou plusieurs groupes. De même un message peut être adressé à un ou plusieurs groupes destinataires.

Un service de communication de groupe simplifie la tâche du concepteur d'applications en lui offrant deux fonctionnalités majeures. D'une part, celui-ci a la possibilité de créer des groupes de processus dont la composition évolue dynamiquement. D'autre part, il peut définir des contraintes concernant l'ordre de livraison des messages de l'application par chaque membre du groupe.

- A tout instant, un processus peut devenir membre d'un groupe en émettant une requête *join* à destination de ce groupe. De même, il peut volontairement cesser d'être membre du groupe en émettant une requête *leave* à destination du groupe. Il peut également cesser d'être membre du groupe à la suite d'une défaillance (panne franche) dont l'occurrence n'est bien évidemment pas prévisible. Tous ces changements de la composition du groupe doivent être observés de manière cohérente par l'ensemble des membres du groupe.
- Les messages destinés à un groupe ne sont pas nécessairement reçus dans le même ordre par chacun des membres du groupe. Ceci est en particulier le cas lorsque les membres d'un même groupe sont situés sur des sites géographiquement distants. Le concepteur d'applications peut imposer que les requêtes soient délivrées à l'application dans le même ordre

par chacun des membres. Cette fonctionnalité permet aux différents membres d'évoluer de façon cohérente.

### 4.3 Utilisations considérées

Le service de communication de groupe est un service générique dont les champs d'application sont vastes.

La tolérance aux défaillances est le domaine d'utilisation le plus couramment cité. Pour pouvoir tolérer des défaillances de type pannes franches, un serveur critique est répliqué sur plusieurs sites géographiquement distants. L'ensemble des copies constitue alors un groupe au sein duquel la panne d'une des copies n'a pas d'incidence sur l'exécution des autres copies. Dans un schéma de réplication dite *active*, toutes les copies exécutent les requêtes émises par les clients. Partant de l'hypothèse que le serveur répliqué a un comportement déterministe, toutes les copies passent nécessairement par la même séquence d'états dès lors que l'on garantit que ces requêtes sont exécutées dans le même ordre par chacun des membres du groupe.

Le travail coopératif assisté par ordinateur (Computer Supported Cooperative Work) connaît depuis quelques années un développement important. L'objectif est de permettre à des groupes d'utilisateurs de collaborer à des buts communs. Il peut s'agir par exemple de l'édition simultanée d'un même document par plusieurs co-auteurs ou du développement d'un projet (de CAO, de génie logiciel) faisant appel à la compétence de différentes équipes, chaque équipe étant elle-même constituée de plusieurs membres. La répartition géographique des différents intervenants des groupes induit naturellement une dimension répartie. Par ailleurs, la collaboration repose généralement sur l'utilisation d'objets partagés et pose inévitablement le problème de la cohérence de ces objets et de leur capacité à résister aux défaillances. Outre ces contraintes de cohérence, certaines applications coopératives procèdent de façon intensive à des changements de la composition des groupes. C'est en particulier le cas des applications s'appuyant sur la métaphore du bâtiment virtuel. Le groupe d'individus réunis au sein d'une même pièce (bureau ou salle de réunion virtuelle) constitue un groupe dont la composition évolue continuellement en fonction des arrivées et des départs. Dans un tel contexte, la gestion automatique des changements de vue qui est réalisée par le service de communication de groupe s'avère alors être un service majeur.

## 5 Logiciels

### 5.1 EDEN : un service de communication de groupe

**Participants** : Emmanuelle Anceaume, Francisco Brasileiro, Fabiola Greve, Michel Hurfin [correspondant], Jean-Pierre Le Narzul, Frédéric Tronel, Yun Wang.

**Mots clés** : applications réparties, résistance aux défaillances, temps réel, communication de groupe, consensus, CORBA, hétérogénéité.

**Résumé** : *Depuis plus de deux ans, un effort important est fourni pour développer dans un environnement CORBA un prototype convaincant qui intègre les*

*solutions que nous proposons pour garantir la sûreté de fonctionnement et le respect de contraintes temporelles. Notre objectif dans les années à venir est d'obtenir un démonstrateur répondant aux exigences mentionnées dans l'appel à propositions rédigé par l'OMG sur le thème de la tolérance aux défaillances dans CORBA.*

**Description de EDEN** La plate-forme EDEN offre des services de communication de groupe (construction d'un ordre total sur les requêtes adressées au groupe, gestion de la composition du groupe et de l'installation des vues). L'originalité de la plate-forme réside d'une part dans le fait que tous les services offerts sont conçus comme des instances d'un service d'accord élémentaire et d'autre part dans son architecture qui repose sur un mécanisme de communication par événements. Le système réparti servant de support à l'exécution du groupe est supposé *asynchrone* : aucune hypothèse n'est faite concernant les temps de transfert des messages ou les vitesses d'exécution des membres du groupe qui peuvent éventuellement s'exécuter sur des machines hétérogènes.

La plate-forme EDEN a été développée en utilisant un langage orienté objet (Java). EDEN est composée de deux modules appelés EVA et ADAM. EVA (EVENt based Architecture) est un framework utilisant le concept de communication par événements pour permettre au programmeur une mise en oeuvre simple et efficace de protocoles de communication spécialisés. Les services de communication de groupe développés à l'aide de EVA sont regroupés au sein du module appelé ADAM.

– EVA

Un protocole de communication de groupe est un exemple typique de protocole de communication spécialisé de haut niveau. La structuration par couche qui est habituellement utilisée pour spécifier les protocoles de communication élémentaires, présente de nombreux défauts lorsqu'il s'agit de concevoir des protocoles de plus haut niveau de la même façon. Pour éviter les inconvénients d'une architecture à couche tout en gardant les avantages d'une décomposition d'un problème en services de base indépendants, nous avons opté pour une architecture fondée sur le concept de notification d'événements. Les interactions entre les entités qui mettent en oeuvre la fonctionnalité désirée sont obtenues via des productions d'événements (données en sortie) par des entités productrices et des consommations d'événements (données en entrée) par des entités consommatrices. Un canal d'événement se charge de diriger les événements produits vers les consommateurs intéressés. Cette approche permet de construire des applications performantes, très modulaires et reconfigurables dont les composants peuvent être exécutés en parallèle. Lors du développement d'EVA, nous avons cherché à utiliser des concepts bien établis (service de notification d'événements, *design patterns*) dont la combinaison permet de bâtir des applications de façon simple et méthodique sans pour autant sacrifier l'efficacité du code généré. Une description détaillée d'EVA et de ses principaux atouts est disponible dans [35].

– ADAM

Le module ADAM offre des services de communication de groupe. Plus précisément, il permet d'une part de gérer dynamiquement la composition d'un groupe (ajout et retrait

de membres) et d'autre part de garantir que toutes les requêtes diffusées à l'intention des membres du groupe seront traitées dans le même ordre. Ces requêtes peuvent avoir été générées par un membre du groupe ou par des objets extérieurs au groupe. L'entrelacement entre les changements de composition du groupe et les livraisons de messages respectent les règles classiques définies dans les systèmes de communication de groupe existants (*view synchrony*). L'originalité de la solution proposée vient du fait que tous ces services sont développés au dessus d'un unique service d'accord générique. Les intérêts de cette approche (modularité, meilleure maîtrise du comportement en période de fort asynchronisme, pré-réglage automatique des délais de garde utilisés, détection simple des pertes de messages, possibilité de prendre en compte des contraintes temporelles) sont décrits aux travers des publications du projet parues durant ces trois dernières années.

## 6 Résultats nouveaux

### 6.1 Détection de propriétés

**Participants :** Emmanuelle Anceaume, Jean-Michel Hélary, Michel Raynal.

**Mots clés :** algorithme réparti, causalité, détection de propriétés, détection décentralisée, état global partagé, conjonction de prédicats locaux.

**Résumé :** *Les horloges logiques (scalaires, vectorielles ou matricielles) sont un mécanisme puissant utilisé par de nombreux algorithmes de contrôle réparti. Nous nous sommes intéressés à étudier leur puissance et leur limitation pour détecter des propriétés d'exécutions réparties.*

Les horloges vectorielles constituent le mécanisme le plus approprié pour capter la causalité entre les événements produits lors de l'exécution d'un calcul réparti. Les implémentations traditionnelles des horloges vectorielles nécessitent d'attacher des vecteurs d'entiers de taille  $n$  aux messages de l'application (où  $n$  est le nombre de processus). Dans ce cadre, nous nous intéressons à la capture des dépendances causales sur un sous-ensemble d'événements que l'application considère comme pertinents. Nous présentons un protocole qui fournit un estampillage correct des événements pertinents sous la contrainte suivante : l'information de contrôle attachée à chaque message est bornée à  $b$  identificateurs d'événements (où  $b$  est une constante prédéfinie comprise entre 1 et  $n$ ) [24].

Toutefois, borner *a priori* cette taille peut conduire à une perte irrémédiable d'information si les événements observables ne sont pas correctement choisis. Le défi consiste à déterminer la taille minimale de la surcharge, en fonction d'un ensemble donné d'événements observables, permettant de capter toutes les dépendances causales. Nous répondons à cette question, qui constituait jusqu'alors un problème ouvert. La taille minimale est déterminée à partir d'un graphe valué, qui représente le calcul réparti et ses événements observables. Bien que la construction complète nécessite la connaissance de tout le calcul, le graphe peut être utilisé en différé pour effectuer l'analyse *a posteriori* d'un calcul [25]. Dans [34] nous présentons une famille de protocoles qui associent à chaque événement pertinent une estampille identifiant

de manière unique ses prédécesseurs immédiats. Cette famille est définie par une condition abstraite qui permet aux messages applicatifs de transmettre des informations de contrôle de taille bornée ( $< n$ ,  $n$  étant le nombre de processus). Cette famille définit une suite de protocoles IPT efficaces.

Nous avons aussi étudié la communication logiquement instantanée (LI) [15, 17, 30]. C'est un mode de communication plus contraignant que la communication causale mais moins que le rendez-vous. En effet lors d'un rendez-vous un événement d'émission s'exécute en même temps que l'événement de réception correspondant. La communication LI assure qu'il existe un estampillage logique (référentiel logique de temps) tel que cette contrainte d'instantanéité est vérifiée.

## 6.2 Systèmes transactionnels dupliqués fondés sur la communication de groupe

**Participants :** Udo Fritzke, Philippe Ingels, Achour Mostéfaoui, Michel Raynal.

**Mots clés :** transaction, duplication, diffusion, tolérance aux défaillances.

**Résumé :** *Nous avons étudié le problème de la diffusion fiable atomique dans un groupe et dans plusieurs groupes. Nous avons ensuite implémenté une primitive de diffusion atomique dans plusieurs groupes où des processus peuvent être défaillants. D'autre part, nous avons étudié l'utilisation d'une telle primitive dans le contrôle de duplication dans les systèmes transactionnels.*

Nous avons étudié la diffusion fiable atomique dans les systèmes sujets aux défaillances. Il a été prouvé que ce problème est équivalent au problème de consensus. Nous avons donc cherché les meilleures implémentations en utilisant le consensus comme brique de base. D'abord dans un système à panne franche [29], ensuite dans les cas où les processus peuvent repartir après défaillance [32, 33].

Nous avons proposé un protocole mettant en œuvre une primitive de *diffusion fiable totalement ordonnée* (ou diffusion atomique) vers plusieurs groupes qui serait moins coûteuse et plus résistante au facteur d'échelle [11]. Ce problème est plus dur à résoudre que le précédent. Sa conception repose sur deux briques de base, en l'occurrence, la *diffusion fiable uniforme* et le *consensus uniforme*. La construction de ce protocole assure deux propriétés principales. La *minimalité* : seuls les processus des groupes destinataires d'un message et le processus émetteur participent au protocole de diffusion (réduction du coût). La seconde propriété est la *localité* : un consensus ne peut concerner que les membres d'un même groupe (résistance au facteur d'échelle). Ces travaux ont été implémentés au dessus de TCP. Par ailleurs, des analyses sur les performances du protocole de consensus et de celui de la diffusion atomique ont été menées sur la maquette réalisée.

Nous avons également étudié l'utilisation des primitives de diffusion fiable uniforme et de diffusion atomique dans le support aux protocoles de contrôle de duplication de données [39]. Contrairement aux modèles de duplication complète trouvés dans la littérature, où une base de données est entièrement dupliquée sur un groupe de processus, nous nous sommes basés sur le

modèle de duplication partielle dans lequel chaque objet est dupliqué de manière indépendante sur un certain groupe de processus. Le protocole décrit assure la “sérialisabilité à une copie” des transactions, et cela grâce aux propriétés de la primitive de diffusion atomique et au verrouillage à deux phases. Par ailleurs, on garantit l’absence d’interblocage de transactions en permettant au protocole d’annuler certaines transactions. On tire profit de deux politiques d’annulation pour aboutir à deux variantes de ce protocole de base, l’une privilégiant les opérations de lecture des transactions et l’autre privilégiant les opérations d’écriture. Le premier protocole procure des exécutions dans lesquelles les transactions de lecture seule ne sont jamais annulées. Le deuxième procure des exécutions dans lesquelles les transactions d’écriture seule ne sont jamais annulées. Ces deux algorithmes garantissent que toutes les transactions ne sont pas annulées. [18] présente un protocole de validation atomique non bloquante qui peut servir comme dernière étape dans les travaux précédents pour traiter le cas des annulations de transactions.

### 6.3 Points de contrôle et retour arrière

**Participants :** Jean-Michel Hélary, Achour Mostéfaoui, Michel Raynal.

**Mots clés :** causalité, point de contrôle, cohérence, tolérance aux défaillances.

**Résumé :** *Nous avons essentiellement affiné les résultats essentiels obtenus en 1996 et 1997, notamment dans l’étude de la cohérence forte et dans la recherche d’un ensemble minimal de Z-chemins dont la suppression suffit à assurer la propriété RDT.*

*Nous nous sommes intéressés au problème de la détermination de points de contrôle globaux cohérents selon les deux approches coordination induite par les communications (communication-induced checkpointing) et non coordonnée.*

**Études dans le cadre de la notion classique de cohérence** Nous nous sommes intéressés à la question essentielle: “l’abstraction d’un calcul obtenue en ne considérant que les points de contrôle locaux et leurs relations de dépendance, est-elle cohérente?”. Nous avons examiné deux réponses à cette question, chacune considérant une notion de cohérence particulière pour l’abstraction définie par les points de contrôle.

Dans le premier cas, l’abstraction est cohérente si tout point de contrôle local appartient à au moins un point de contrôle global cohérent. Nous avons caractérisé cette propriété de cohérence par une propriété sur l’estampillage des points de contrôle par des entiers, montré comment cet estampillage permet d’obtenir rapidement “en ligne” des points de contrôle globaux cohérents, et proposé un protocole générique, coordonné par les communications. Ce protocole recouvre de nombreux protocoles connus [13].

Nous avons proposé un protocole non coordonné de prises de points de contrôle locaux. Dans une stratégie non-coordonnée, les processus prennent des points de contrôle locaux de manière indépendante. Le problème est alors le suivant :étant donné un sous-ensemble de points de contrôle locaux (au plus un par processus), la théorie de Netzer et Xu permet de répondre à la question *existe-t-il un point de contrôle global cohérent qui contient ce sous-ensemble?* par une caractérisation théorique. Le protocole proposé est capable de répondre *en ligne* à cette

question; de plus, si la réponse est positive, il fournit le point de contrôle global cohérent *le plus récent* englobant le sous-ensemble donné, sinon il fournit le “prochain” point de contrôle global cohérent [8].

La seconde notion de cohérence, plus forte que la précédente, est la propriété RDT (pour *Rollback-Dependency Trackability*) introduite par Wang. Dans le cadre de l’approche *coordination induite par les communications*, nous nous sommes intéressés à la recherche de conditions optimales sous lesquelles un processus sera obligé de prendre un point de contrôle forcé. Par là nous entendons la recherche de conditions telles que le nombre de points de contrôle forcés serait le plus faible possible. Dans une approche “théorique”, il s’agit de définir un ensemble de Z-chemins dont la suppression (par prise de points de contrôles locaux) suffit à assurer la propriété RDT. Nous avons donc caractérisé un ensemble *minimal* de Z-chemins dont la suppression assure la propriété RDT et proposé un protocole induit par les communications qui réalise “au mieux” cette caractérisation optimale [10]. Nous avons de plus établi un résultat d’impossibilité sur l’utilisation des horloges scalaires pour garantir la propriété RDT [9].

## 6.4 Les problèmes d’accord dans les systèmes répartis

**Participants :** Emmanuelle Anceaume, Francisco Brasileiro, Udo Fritzke, Fabiola Greve, Jean-Michel Hélary, Michel Hurfin, Philippe Ingels, Jean-Pierre Le Narzul, Achour Mostéfaoui, Michel Raynal, Frédéric Tronel, Yun Wang.

**Mots clés :** algorithme réparti, résistance aux défaillances, duplication, communication de groupe, temps réel, consensus.

**Résumé :** *Les travaux réalisés au cours de l’année 2000 sur les problèmes d’accord ont porté sur trois aspects. Le premier concerne la solvabilité du consensus (nouvelles familles de détecteurs de défaillances, introduction des probabilités et enfin utilisation de conditions sur l’ensemble des valeurs proposées). Le second aspect concerne l’utilisation du problème du consensus comme brique de base pour résoudre des problèmes d’accord. Enfin, le troisième aspect concerne l’élargissement du champ d’application de nos résultats en considérant les fautes byzantines et les contraintes temporelles. Les différents résultats obtenus ont, en partie, été intégrés à la plate-forme EDEN que nous développons.*

Si aucune hypothèse supplémentaire n’est adoptée, il est prouvé que le problème d’accord élémentaire n’admet pas de solution déterministe dans un système réparti asynchrone dès lors qu’un processus au moins est susceptible de connaître une défaillance de type panne franche. Notre axe de recherche principal est donc naturellement consacré à l’étude des solutions permettant de contourner ce résultat d’impossibilité principalement les solutions déterministes fondées sur les détecteurs de défaillances.

### – Solvabilité du consensus

Afin de contourner le résultat d’impossibilité nous avons suivi trois directions. La première direction, initiée par Chandra et Toueg, consiste à enrichir le système sous-jacent

par des détecteurs de défaillances. Dans [16], nous nous intéressons à de nouvelles classes de détecteurs de défaillances dont la portée est restreinte. Plus précisément, on suppose que les propriétés de complétude (respectivement, d'exactitude) telles que définies classiquement par Chandra, Hadzilacos et Toueg ne portent pas sur l'ensemble des processus corrects mais seulement sur  $k$  (respectivement,  $k'$ ) processus corrects parmi eux. On donne les conditions sous lesquelles ces détecteurs restreints sont équivalents à leurs homologues de portée maximale. Dans [28], les restrictions ne portent plus seulement sur les processus corrects mais sur tout processus. Les détecteurs de défaillance sont donc encore plus simples à mettre en œuvre. [46] présente une introduction à des détecteurs de défaillances pour résoudre le problème de la diffusion fiable (plus simple à résoudre que le consensus). Ces derniers peuvent être mis en œuvre moyennant une majorité de processus corrects.

Une autre façon de contourner le résultat d'impossibilité consiste à affaiblir le problème (la terminaison est seulement probabiliste). Ben-Or a été le premier à trouver une solution probabiliste au problème du consensus binaire. Tout comme Aguilera et Toueg, nous avons exploré dans [26] une approche hybride pour résoudre le problème du consensus. La solution proposée combine dans la même structure algorithmique deux sources possibles de terminaison. La première, déterministe, est fondée sur des détecteurs de défaillances. La seconde, probabiliste, est fondée sur le tirage aléatoire. A cause de sa composante probabiliste, le protocole hybride ne résout que le consensus binaire. Nous avons donc proposé un protocole simple qui résout le consensus multi-valué en utilisant un protocole de consensus binaire comme brique de base [14]. Afin d'éviter cet empilement protocolaire, nous avons proposé un protocole de consensus multi-valué purement probabiliste [38].

Nous avons introduit une nouvelle approche dans la résolution du problème du consensus dans les systèmes répartis asynchrones [43]. Un vecteur d'entrée au problème du consensus est le vecteur constitué des  $n$  valeurs proposées par les  $n$  processus. Si  $V$  est l'ensemble des valeurs qui peuvent être proposées alors il existe  $|V|^n$  vecteurs d'entrée possibles au consensus. Une manière de circonvenir le résultat d'impossibilité consiste à n'autoriser qu'un sous-ensemble  $S$  de ces vecteurs. Le premier résultat que nous proposons est un protocole de consensus générique, paramétré par  $S$ , pour le modèle à mémoire partagée. Ce protocole garantit la propriété d'accord du consensus quel que soit le vecteur d'entrée. La propriété de terminaison est satisfaite soit en l'absence de défaillances soit lorsque le vecteur d'entrée appartient au sous-ensemble  $S$ . Nous avons donné une condition nécessaire et suffisante que doit vérifier un sous-ensemble de vecteurs d'entrée pour résoudre le problème du consensus. Ces travaux effectués pour le modèle à mémoire partagée ont été étendus au modèle de communication par messages.

Les pannes franches ne constituent qu'un cas particulier de comportement défaillant. Aussi, concevoir des protocoles répartis capables de résister à un comportement erroné arbitraire des processus constitue un réel défi, compte-tenu de la possibilité d'attaques malveillantes ou d'erreurs logicielles imprévues. Par conséquent, être capable de construire un protocole masquant les défaillances arbitraires au-dessus de protocoles tolérant les pannes franches constituerait une avancée majeure dans le domaine de l'ingénierie logi-

cielle. Nous avons exploré une approche méthodologique générale permettant de telles constructions [20]. Elle s'applique au cas des protocoles dont tous les processus ont le même programme et échangent régulièrement des messages. Cette approche, modulaire, est basée sur l'encapsulation de la détection de défaillances de tout type dans des modules spécifiques. Elle peut constituer un point de départ pour la conception d'outils de transformation automatique. Cette méthodologie a été appliquée au problème du consensus [19].

#### – Solutions au problème du consensus

Nous nous sommes principalement intéressés aux solutions déterministes s'appuyant sur des détecteurs de défaillances appartenant aux classes  $\mathcal{S}$  et  $\diamond\mathcal{S}$ . La classe des détecteurs de défaillances forts (dénotee  $\mathcal{S}$ ) regroupe tous les détecteurs de défaillances qui ont pour propriétés de suspecter les processus défaillants (propriété de complétude) mais de ne pas suspecter au moins un des processus non défaillants (propriété d'exactitude). La classe des détecteurs de défaillances faibles (dénotee  $\diamond\mathcal{S}$ ) regroupe quant à elle tous les détecteurs de défaillances qui vérifient la même propriété de complétude mais ne satisfont la propriété d'exactitude qu'après un laps de temps indéterminé. De fait, les détecteurs de défaillances appartenant à ces deux classes sont intrinsèquement non fiables puisqu'ils peuvent suspecter arbitrairement des processus corrects. Bien que faibles, les propriétés d'exactitude caractérisant ces deux classes sont suffisantes pour permettre de résoudre le problème du consensus si au moins un processus est correct (classe  $\mathcal{S}$ ) ou si au moins une majorité de processus est correcte (classe  $\diamond\mathcal{S}$ ).

Plusieurs protocoles de consensus fondés sur la classe  $\mathcal{S}$  ont été conçus par le passé, dont un très simple que nous avons proposé récemment [27]. Certains d'entre-eux requièrent systématiquement  $n$  étapes de calcul ( $n$  étant le nombre de processus), chaque étape nécessitant l'émission de  $n$  ou  $n^2$  messages. D'autres protocoles permettent une prise de décision au plus tôt (dans ce cas, le nombre d'étapes de calcul dépend du nombre de défaillances détectées) mais requièrent  $n^2$  messages à chaque étape de calcul. Une de nos contributions [21, 22] a consisté à définir un protocole de consensus fondé sur la classe  $\mathcal{S}$  qui permet des décisions au plus tôt. Chaque étape de calcul nécessite  $3(n-1)$  messages. De fait, le protocole que nous proposons est particulièrement efficace en temps et en nombre de messages. Ce protocole a été adapté aux détecteurs de défaillances de la classe  $\diamond\mathcal{S}$  et généralisé afin de permettre, lorsque cela est nécessaire, une nouvelle réduction du nombre d'étapes de calcul au prix bien évidemment d'un accroissement du nombre de messages échangés par étape. Le compromis qui résulte du choix d'un schéma d'échange de messages peut être modifié à chaque pas d'itération. Ainsi le nombre de messages échangés par pas de communication peut varier de  $O(n)$  (schéma centralisé) à  $O(n^2)$  (schéma totalement décentralisé) [41].

En plus de ce premier degré de liberté, nous proposons dans [42], un protocole générique qui peut s'adapter aisément aux deux classes de détecteurs de défaillances  $\diamond\mathcal{S}$  et  $\mathcal{S}$ . Différentes instantiations du protocole général permettent de retrouver des protocoles existants (par exemple, celui de Chandra et Toueg) ainsi que de nouveaux protocoles. Cette approche, également intéressante par son aspect méthodologique, permet de com-

prendre les fondements conceptuels de cette famille de protocoles. Ceci est rendu possible par la caractérisation de deux ensembles de processus qui, pour un pas d'itération donné, assurent respectivement les propriétés de vivacité et de sûreté. Cette approche permet également de montrer que, contrairement à une idée admise, les estampilles utilisées dans le protocole de Chandra et Toueg peuvent être bornées.

Un des défauts des protocoles fondés sur la classe  $\diamond\mathcal{S}$  réside dans le fait qu'à une étape donnée seule la valeur du coordonnateur peut être décidée. Si nous considérons le cas favorable où tous les processus proposent la même valeur, il est malgré tout possible de ne pas décider avant longtemps si la propriété d'exactitude met du temps à se vérifier. Nous avons donc défini un protocole [36] qui permet aux processus de décider en une étape de communication dans des circonstances favorables. Celles-ci apparaissent lorsque  $(n - f)$  processus proposent la même valeur (où  $f$  est le nombre maximal défaillances). Le protocole requiert la condition  $f < n/3$ . Il est montré que, dans le cas général, cette condition est nécessaire. Cette condition peut être affaiblie à  $f < n/2$  au prix d'une étape de communication supplémentaire.

#### – Autres problèmes d'accord

Le consensus est un problème d'accord élémentaire dont la définition simple est unanimement reconnue par tous. Il s'agit d'un problème purement théorique qui permet d'abstraire différents problèmes d'accord. La simplicité de ce problème permet, lors de son étude, de se focaliser uniquement sur le résultat d'impossibilité qui est associé à tout problème d'accord. En pratique, ceci ne signifie pas pour autant qu'un service de consensus est le dénominateur commun idéal à partir duquel on peut toujours dériver des solutions efficaces à des problèmes d'accord plus complexes. Diverses variantes au problème du consensus ont donc été explorées. Le but poursuivi est de concevoir une brique de base parfaitement adaptée au problème d'accord que l'on désire résoudre.

Parmi les variantes possibles au problème du consensus, nous avons étudié le problème du calcul d'une donnée globale cohérente (un vecteur comportant une entrée par processus). L'importance du problème nous a conduit à rechercher un protocole spécifique, efficace en temps, utilisable uniquement si on dispose de détecteurs de défaillances parfaits [12, 23]. Ce protocole permet des décisions au plus tôt. Soit  $t$  le nombre maximal de pannes pouvant se produire et soit  $f$  le nombre de pannes effectives ( $f < t$ ). Dans le pire des cas, le protocole termine après  $\min(2f + 2, t + 1)$  tentatives. De plus, ce protocole ne nécessite pas d'échange d'information sur la perception locale des pannes. Ce travail présenté à la conférence ICDCS'00 (International conference on Distributed Computing Systems) s'est vu attribué le "Best Paper Award".

Le problème du " $k$ -Set Agreement" généralise le problème du consensus en ce sens qu'il ne limite pas l'ensemble des valeurs de décision à une seule mais à  $k$  valeurs, il est de ce fait moins dur à résoudre. Dans un système asynchrone composé de  $n$  processus et où au plus  $f$  d'entre eux peuvent être défaillants, le problème du " $k$ -Set Agreement" peut être facilement résolu si  $f < k$ . Il a par ailleurs été démontré qu'aucune solution déterministe n'existe lorsque  $f \geq k$ .

Nous avons montré dans [28] que les détecteurs de défaillances de Chandra et Toueg

ne sont pas minimaux pour résoudre le “ $k$ -Set Agreement” nous avons donc étudié une approche fondée sur les détecteurs de défaillances à portée réduite. Dans [45] nous présentons une approche probabiliste pour résoudre ce problème. Le protocole proposé ne nécessite pas la connaissance *a priori* de l’ensemble des valeurs proposées. Dans les deux cas, nous avons mis en évidence un aspect contre-intuitif du problème : plus “ $k$ ” augmente moins on tolère de défaillances. Cependant, on gagne sur l’aspect impossibilité puisque dans la solution déterministe on a besoin de détecteurs de défaillances plus faibles. Tout comme dans le cas du consensus, on a présenté dans [44] un protocole fondé sur une restriction sur la composition de l’ensemble des valeurs proposées qui permet de résoudre ce problème quel que soit  $f$ . En l’occurrence, nous avons montré que s’il existe parmi toutes les valeurs proposées un sous-ensemble de  $k$  valeurs proposées par plus de  $(kn + f)/(k + 1)$  processus alors le problème du “ $k$ -Set Agreement” a une solution.

#### – Prise en compte des contraintes temporelles

Notre objectif est de maîtriser la terminaison des protocoles exécutés en permettant de gérer au mieux une durée maximale d’exécution.

Un protocole fondé sur l’ordre causal garantit que si deux messages sont causalement liés et ont les mêmes destinataires, alors ces messages sont remis (à l’application) selon leur ordre d’émission. Cette garantie permet de simplifier sensiblement le développement d’applications distribuées. Pour empêcher la violation de l’ordre causal, deux techniques existent : soit on force un message à attendre les messages “(en retard)” dont il dépend, soit les messages “en retard” sont annulés. Il est clair que pour un système temps réel, la première approche n’est pas satisfaisante. En effet dès qu’un message dépasse son échéance, tous les messages qui en dépendent vont forcément dépasser leur échéance. Nous proposons dans [31] une abstraction originale de l’ordre causal fondé sur les échéances des messages. Deux implémentations sont proposées dans le contexte des communications de groupe (diffusion générale et diffusion sélective).

#### – Gestion de la composition d’un groupe

Le concept de groupe est de plus en plus utilisé pour concevoir des couches logicielles permettant de supporter l’exécution d’applications réparties fiables. Ce paradigme recouvre deux services de base, à savoir, un service de gestion de la composition du groupe et un service de communication de groupe. Un groupe est un ensemble de processus coopérant à la réalisation d’une tâche commune. Étant donné que de nouveaux processus peuvent souhaiter rejoindre le groupe, que des membres du groupe peuvent souhaiter le quitter ou connaître des défaillances de type panne franche, la composition du groupe doit pouvoir évoluer de façon dynamique. L’ensemble des processus qui constituent le groupe à un instant donné est appelé la vue courante du groupe.

Dans [40], nous nous intéressons aux problèmes de la spécification et du développement d’un service de gestion de la composition du groupe ne tolérant qu’une partition primaire (une seule vue à un instant donné). Nous proposons tout d’abord une spécification du problème. Ensuite, nous présentons un protocole qui satisfait cette spécification dans un système réparti asynchrone équipé de détecteurs de défaillances. Le service de gestion

de la composition du groupe via une partition primaire est obtenu en spécialisant de manière appropriée un protocole d'accord générique. Pour cela, nous avons défini un cadre général dans lequel la définition d'un problème d'accord particulier est fixée par le biais de 6 paramètres. Le protocole proposé par Chandra et Toueg pour résoudre le problème du consensus à l'aide de détecteurs de défaillance appartenant à la classe  $\diamond\mathcal{S}$  constitue l'ossature du protocole générique proposé. Les extensions offertes par ce nouveau service d'accord permettent notamment aux membres d'un groupe de processus de décider unanimement sur une collection de valeurs proposées (plutôt que sur une seule des valeurs proposées) tout en autorisant les processus à effectuer de multiples propositions (plutôt qu'une seule au démarrage du protocole d'accord).

## 6.5 Aide à la conception des applications multimedia

**Participants :** Emmanuelle Anceaume, Erwan Demairy.

**Mots clés :** application multimedia, architecture qualité de services, temps réel, allocation de ressources.

**Résumé :** *Les travaux réalisés ont tout d'abord porté sur l'aspect "statique" d'une application multimedia puis sur son aspect "dynamique". Plus précisément, nous avons conçu un algorithme permettant de vérifier statiquement une application vis-à-vis de la cohérence des données échangées et du respect des contraintes temporelles. Nous nous sommes ensuite intéressés à la répartition dynamique des ressources entre les applications qui s'exécutent sur une architecture qualité de service donnée.*

Nos travaux s'inscrivent dans le cadre de l'aide à la conception des applications multimedia. Nous nous intéressons plus particulièrement aux applications multimedia capables d'adaptation entre différents niveaux de qualité. Afin de faciliter le travail des concepteurs d'applications multimedia, nous définissons tout d'abord un algorithme qui permet de vérifier statiquement une application vis-à-vis de la cohérence des données échangées et du respect des contraintes temporelles. L'application est décrite par son concepteur sous la forme de composants et de connecteurs, habituelle dans le domaine des architectures logicielles. On associe à ces éléments les configurations qu'ils supportent, c'est-à-dire la qualité qu'ils peuvent fournir ou accepter. Chacune de ces configurations impose une exigence en ressources pour le système sous-jacent. Ces exigences sont obtenues en mesurant la consommation des applications, par le biais de moniteurs. À ce stade, nous pouvons vérifier que l'association des composants et connecteurs constitue bien un ensemble cohérent vis-à-vis des qualités comprises et fournies par les éléments constituant l'application. D'autre part, des contraintes temporelles probabilistes sont placées sur les flux de l'application. Nous définissons une méthode qui permet de déterminer si les contraintes peuvent être respectées, en utilisant les garanties données par les éléments constituant l'application.

Ce premier algorithme fonctionne en considérant que les ressources nécessaires pour garantir le bon fonctionnement des éléments de l'application sont disponibles. Cette hypothèse n'étant

en général pas satisfaite, nous nous intéressons dans un deuxième temps, à la répartition des ressources entre applications. Nous supposons que ces applications s'exécutent sur une architecture qualité de service donnée. Cet algorithme est utilisé par l'environnement pour piloter une AQDS au moment de l'admission ou du retrait d'une application en permettant une redistribution des ressources entre les différentes applications. L'environnement agit comme une sur-couche d'une architecture de qualité de service et permet d'obtenir une qualité "optimale"; le critère d'optimalité étant précisé ultérieurement.

## 7 Contrats industriels (nationaux, européens et internationaux)

### 7.1 Contrat FT R&D sur les contraintes non-fonctionnelles

**Participants :** Fabiola Greve, Michel Hurfin, Eric Mourgaya, Michel Raynal, Frédéric Tronel.

**Mots clés :** sûreté de fonctionnement, temps réel, CORBA.

**Résumé :** *Suite aux Consultations Thématiques Informelles lancées par le CNET en mars 1997, un contrat a été établi entre le CNET et l'équipe ADP pour les années 1998 à 2001 (convention 1 98 C 172 00 31401 012). Ce contrat porte sur l'étude et le développement d'un ensemble de services garantissant, dans un environnement asynchrone composé de calculateurs hétérogènes, deux types de contraintes non-fonctionnelles à savoir la sûreté de fonctionnement et le temps réel. Ce contrat doit déboucher sur la définition d'un ensemble de protocoles et la réalisation d'un prototype qui sera développé dans un environnement propre au CNET (Distributed Processing Environment Jonathan).*

### 7.2 Contrat RÉUTEL 2000

**Participants :** Francisco Brasileiro, Udo Fritzke, Fabiola Greve, Michel Hurfin, Philippe Ingels, Jean-Pierre Le Narzul, Frédéric Tronel, Yun Wang.

**Mots clés :** sûreté de fonctionnement, temps réel, CORBA.

**Résumé :** *L'équipe ADP est responsable de l'une des quatre actions définies dans le contrat Réutel 2000, qui a été signé en 1997. Dans le cadre de ce contrat ALCATEL ALSTHOM RECHERCHE souhaite bénéficier des travaux de recherche que nous effectuons dans les domaines de la reprise après défaillance et de la communication de groupe. Les travaux réalisés ont pour objectif d'enrichir les services offerts dans un environnement CORBA en proposant des extensions logicielles qui permettront aux programmeurs de définir des services tolérant les défaillances (duplication active). L'approche proposée consiste à concevoir des services nécessaires à la gestion d'un groupe de serveurs dupliqués (e.g., gestion de la composition du*

*groupe, diffusion atomique) en considérant ceux-ci comme des variantes d'un problème d'accord élémentaire : le consensus. Notre travail consiste donc à développer un prototype en y intégrant l'ensemble des innovations proposées tout en respectant la spécification d'un service CORBA de tolérance aux défaillances qui a été adoptée par l'OMG au début de l'année 2000.*

## 8 Actions régionales, nationales et internationales

### 8.1 Actions nationales

#### 8.1.1 GDR ARP (Architecture, Réseaux et Parallélisme) du CNRS

**Participants :** Achour Mostéfaoui, Michel Raynal.

Le thème Réseaux et Systèmes de ce GDR est coordonnée par Michel Diaz (LAAS). Le sous-groupe "Algorithmes pour les Systèmes Répartis" est animé par Joffroy Beauquier (LRI) et Michel Raynal (IRISA).

### 8.2 Actions européennes

#### 8.2.1 Cabernet

**Participants :** Michel Hurfin, Michel Raynal.

Le projet ADP participe au projet Esprit Cabernet (*Computer Architecture for Basic European Research, Network of Excellence*) en collaboration avec le projet SOLIDOR de l'IRISA et l'action REFLECS de l'INRIA-Rocquencourt depuis sa création en 1991.

#### 8.2.2 Italie

**Participants :** Emmanuelle Anceaume, Jean-Michel Hélary, Michel Raynal.

Le projet de coopération avec l'Université de Rome La Sapienza sur le thème des points de contrôle et du retour arrière qui avait été retenu pour financement dans le cadre des programmes d'accord CNRS-CNR 1998-1999 a été reconduit pour 2000-2001.

### 8.3 Actions internationales

#### 8.3.1 Amériques

##### – USA (Géorgie)

**Participants :** Achour Mostéfaoui, Michel Raynal.

Un projet de coopération avec M. Ahamad du Georgia Institute of Technology, sur le thème de la cohérence des mémoires réparties, a été accepté par la NSF et l'INRIA pour un financement conjoint. Ce projet a débuté en 1998.

– **Mexique**

**Participants** : Emmanuelle Anceaume, Achour Mostéfaoui, Michel Raynal.

Un projet de coopération avec S. Rajsbaum et M. Romero de l'UNAM, Université nationale de Mexico, sur le thème des problèmes d'accord dans les systèmes asynchrones, a été retenu pour financement dans le cadre des programmes d'accords CNRS-CONACYT 2000-2001.

– **Brésil (Université de Bahia)**

**Participants** : Fabiola Greve, Michel Hurfin, Michel Raynal, Frédéric Tronel.

Un projet de coopération avec R. Macedo de l'Université Fédérale de Bahia (Salvador) a été retenu dans le cadre du programme de coopération INRIA/CNPq. L'objectif de ce projet d'une durée de deux ans (2000/2001) est de concevoir et de développer un prototype commun offrant des services de communication de groupe. Côté français, ce projet est géré par Michel Hurfin.

– **Brésil (Université de Santa Catarina)**

**Participants** : Michel Hurfin, Michel Raynal.

Pour l'année 2000, un projet de coopération avec J.M. Farines de l'Université Fédérale de Santa Catarina (Florianopolis) a été retenu dans le cadre des accords CNRS/CNPq. Ce projet porte sur l'étude des contraintes non fonctionnelles dans les environnements CORBA. Côté français, ce projet est géré par Michel Raynal.

– **Chine**

**Participants** : Michel Hurfin, Jean-Pierre Le Narzul, Achour Mostefaoui, Michel Raynal, Frédéric Tronel.

Un projet de coopération intitulé "Fault Tolérant Corba", dont Michel Hurfin est le coordinateur coté Français, a été sélectionné par le LIAMA (Laboratoire franco-chinois d'informatique, d'automatique et de mathématiques appliquées) pour bénéficier d'un financement durant une année 2000. L'objectif est d'intégrer les mécanismes de communication de groupe que nous développons au sein de l'ORB développé par l'équipe chinoise en respectant les spécifications récemment adoptées par l'OMG.

## 9 Diffusion de résultats

### 9.1 Actions d'enseignement

A l'exception de deux chercheurs à temps plein, tous les membres permanents du projet ADP sont enseignants-chercheurs. A ce titre, ils participent à bon nombre d'enseignements de

second et de troisième cycles. Par ailleurs J.M. Hélarly est responsable des enseignements du diplôme d'ingénieur en informatique de l'université de Rennes 1 (DIIC) et A. Mostefaoui est responsable de la filière Langages et Systèmes Informatiques du DIIC.

Le projet a assuré plusieurs enseignements liés à ses thèmes de recherche :

- dans le DEA d'informatique de Rennes, M. Raynal dispense un cours sur les algorithmes et systèmes répartis (21 heures) ;
- dans le DIIC 3<sup>e</sup> année (filières LSI et ARC), J.M. Hélarly et M. Raynal dispensent un cours sur l'algorithmique répartie (30 heures) ;
- à l'ENSTBr (antenne de Rennes), un cours sur l'algorithmique répartie est assuré par M. Raynal. M. Hurfin et M. Raynal dispensent également un cours sur la tolérance aux défaillances à l'ENSTBr de Brest (respectivement de 3 et 6 heures chacun) ;
- dans le DESS ISA, M. Raynal donne un cours sur les systèmes transactionnels, la cohérence et la sécurité des données.
- J.-M. Hélarly a dispensé un cours sur l'algorithmique répartie (15 heures) dans le cadre du DEA MILS (Modélisation et Ingénierie du Logiciel Scientifique) de l'université libanaise de Beyrouth.
- M. Hurfin a dispensé un cours sur la tolérance aux défaillances (15 heures) aux étudiants de Southeast University (Nankin) en Chine. Une présentation sur ce même sujet a été faite à Tsinghua University (Pékin), et au Beijing Institute of Radio Measurment (Pékin).

Un stagiaire de DEA a effectué son stage au sein de l'équipe ADP. Il s'agit de Eric Mourgaya qui a modélisé par des chaînes de Markov trois protocoles de consensus représentatifs de l'état de l'art.

## 9.2 Présentation de travaux

Les membres de l'équipe ont participé à diverses conférences et *workshops* (se reporter à la bibliographie pour en avoir la liste).

M. Raynal a présenté les travaux de l'équipe et fait des séminaires dans les universités et centres de recherche suivants :

- Université du Connecticut (USA), jan. 2000,
- Georgia Institute of Technology, (Atlanta, Géorgie), mars 2000,
- Université La Sapienza de Rome (Rome, Italie), juin 2000.

### 9.3 Animation de la communauté scientifique

M. Raynal est membre du comité de lecture des revues suivantes *Foundations of Computer and Decision Sciences*, *Journal of Computer Systems Science and Engineering* et *Journal of Distributed Systems Engineering*.

M. Raynal a été (en 2000) ou sera (en 2001) membre des comités de programme des conférences suivantes :

- *20th IEEE Int. conf. on Distributed Computing Systems*.  
Taiwan, avril 2000. (actes publiés par IEEE).
- Co-chair (avec M. Takizawa, Tokyo Denki University) du *Workshop on Group Computation and Communications*.  
Taiwan, avril 2000. (Actes publiés par IEEE). (Ce workshop est une manifestation satellite du *20th IEEE Int. conf. on Distributed Computing Systems*).
- *3rd IEEE Int. Symposium on Object-Oriented Real-time Distributed Computing (ISORC'2K)*.  
Newport Beach, California, mars 2000. (Actes publiés par IEEE).
- *Int. Conf. on Dependable Systems and Networks (IEEE Symposium on Fault-Tolerant Systems FTCS'30 and DCCA'8)*.  
New York, juin 2000. (Actes publiés par IEEE)
- *Fifth International Symposium on Autonomous Decentralized Systems*.  
Dallas, Texas, mars 2001.
- Publicity Chair *21th IEEE Int. conf. on Distributed Computing Systems*.  
Phoenix, Arizona, avril 2001, (Actes publiés par IEEE)
- *20th ACM Symposium on Principles of Distributed Computing (PODC'01)*.  
Newport, Road Island, août 2001, (Actes publiés par ACM).
- co-président et co-organisateur (avec Luis Rodrigues, Lisboa, PT) du *Int. IEEE Workshop on Applied Reliable Group Communication* (in conjunction with the 21th IEEE Int. conf. on Distributed Computing Systems.  
Phoenix, Arizona, avril 2001. (Actes publiés par IEEE).
- *2d Int. Conf. on Dependable Systems and Networks (IEEE Symposium on Fault-Tolerant Systems FTCS'31 and DCCA'9)*.  
Goteborg, Juillet 2001. (Actes publiés par IEEE).
- *20th IEEE Symposium on Reliable Distributed Systems (SRDS'01)*.  
Actes publiés par IEEE).
- Program Chair, *22th IEEE Int. conf. on Distributed Computing Systems*.  
Vienne, Autriche, juillet 2002. (Actes publiés par IEEE).

Michel Hurfin est membre du comité de programme de ICDCS-2001 (21st International Conference on Distributed Computing Systems) qui se déroulera à Phoenix en Arizona en avril 2001.

A. Mostéfaoui a été élu au Conseil d'Administration de l'Association SPECIF des enseignants et des chercheurs en Informatique.

E. Anceaume et M. Hurfin ont mis en place, organisé et animent en coopération avec L. Toutain de l'ENSTbr le Séminaire Réseaux et Systèmes (<http://www.irisa.fr/adp/srs/srs.html>).

## 9.4 Coopération et échanges internationaux

Suite à la visite de Amr El Abbadi, Professeur à l'Université Santa Barbara (Californie) au cours des mois de juillet et août 1999, une proposition de coopération NSF/INRIA a été déposée en juillet 2000. Le thème de cette proposition concerne les problèmes liés à la mise à jour des "warehouse".

# 10 Bibliographie

## Ouvrages et articles de référence de l'équipe

- [1] O. BABAOĞLU, E. FROMENTIN, M. RAYNAL, «A Unified Framework for Expressing and Detecting Run-Time Properties of Distributed Computations», *Journal of Systems and Software, Numéro spécial sur Software Engineering for Distributed Computing* 33, 3, juin 1996, p. 287–298.
- [2] J. BRZEZINSKI, J. M. HÉLARY, M. RAYNAL, M. SINGHAL, «Deadlock models and a general algorithm for distributed deadlock detection», *Journal of Parallel and Distributed Computing* 31, 2, 1995, p. 112–125.
- [3] J. M. HÉLARY, A. MOSTEFAOUI, M. RAYNAL, «A general scheme for token and tree based distributed mutual exclusion algorithm», *IEEE Transactions on Parallel and Distributed Systems* 5, 11, November 1994, p. 1185–1196.
- [4] M. HURFIN, N. PLOUZEAU, M. RAYNAL, «Detecting Atomic Sequences of Predicates in Distributed Computations», in : *Proc. of the ACM Conference on Parallel and Distributed Debugging*, p. 32–42, San Diego, Californie, May 1993. Reprinted in SIGPLAN Notices, vol. 28,12, December 1993.
- [5] M. RAYNAL, A. SCHIPER, S. TOUEG, «The Causal Ordering Abstraction and a Simple Way to implement it», *Information Processing Letters* 39, September 1991, p. 343–351.
- [6] M. RAYNAL, M. SINGHAL, «Logical Time: Capturing Causality in Distributed Systems», *IEEE Computer* 29, 2, février 1996, p. 49–57.
- [7] A. SCHIPER, M. RAYNAL, «From group communication to transactions in distributed systems», *Communications of the ACM* 39, 4, avril 1996, p. 84–90.

## Articles et chapitres de livre

- [8] R. BALDONI, G. CIOFFI, J. HELARY, M. RAYNAL, «Direct Dependency-Based Determination of Consistent Global Checkpoints», *International Journal of Computer and Systems Science and Engineering* 15, 5, Sep 2000.
- [9] R. BALDONI, J. HELARY, A. MOSTEFAOUI, M. RAYNAL, «Impossibility of Scalar Clocked-Based Communication-Induced Checkpointing Protocols Ensuring the RDT Property», *Information Processing Letters*, 2000, To appear.
- [10] R. BALDONI, J. HELARY, M. RAYNAL, «Rollback-dependency trackability: a minimal characterization and its protocol», *Information and Computation*, 2001, To appear.
- [11] U. FRITZKE, P. INGELS, A. MOSTEFAOUI, M. RAYNAL, «Consensus-Based Fault-Tolerant Total Order Multicast», *IEEE Transactions on Parallel and Distributed Systems*, 2000, To appear.
- [12] J. HELARY, M. HURFIN, A. MOSTEFAOUI, M. RAYNAL, F. TRONEL, «Computing Global Functions in Asynchronous Distributed Systems with Perfect Failure Detectors», *IEEE Transactions on Parallel and Distributed Systems* 11, 9, Sep 2000, p. 897–909.
- [13] J. HELARY, A. MOSTEFAOUI, R. NETZER, M. RAYNAL, «Communication-based prevention of useless checkpoints in distributed computations», *Distributed Computing* 13, 1, Jan 2000, p. 29–43.
- [14] A. MOSTEFAOUI, M. RAYNAL, F. TRONEL, «From Binary Consensus to Multivalued Consensus in Asynchronous Message-passing Systems», *Information Processing Letters* 73, 5, 2000, p. 207–212.
- [15] A. MOSTEFAOUI, M. RAYNAL, P. VERISSIMO, «The Logically Instantaneous Communication Mode: a Communication Abstraction», *Future Generation Computing Systems* 17, 7, Apr 2001, To appear.
- [16] M. RAYNAL, F. TRONEL, «Restricted Failure Detectors: Definition and Reduction Protocols», *Information Processing Letters* 72, 3, 2000, p. 91–97.
- [17] M. RAYNAL, *Encyclopedia of Distributed Computing*, Kluwer Academic Press, 2000, ch. Synchronization.
- [18] M. RAYNAL, «Non-Blocking Atomic Commitment in Distributed Systems: A tutorial based on a generic protocol», *Journal of Computer Systems Science and Engineering* 15, 2, 2000, p. 77–86.

## Communications à des congrès, colloques, etc.

- [19] R. BALDONI, J. HELARY, M. RAYNAL, L. TANGUY, «Consensus in Byzantine Asynchronous Systems», in: *Proc of the 7th Int. Colloquium on Structural Information and Communication Complexity (SIROCCO-00)*, p. 1–15, L'Aquila, Jun 2000.
- [20] R. BALDONI, J. HELARY, M. RAYNAL, «From Crash Fault-Tolerance to Arbitrary-Fault-Tolerance: Towards a Modular Approach», in: *Proc. of the Int. Conference on Dependable Systems and Networks (FTCS/DCCA)*, IEEE, p. 273–282, New York, NY, Jun 2000.
- [21] F. GREVE, M. HURFIN, R. MACEDO, M. RAYNAL, «Consensus Based on Strong Failure Detectors: Time and Message Efficient Protocol», in: *Proc. of the Int. Workshop on Fault-Tolerant Parallel and Distributed Systems, LNCS*, 1800, Springer Verlag, p. 1258–1267, Cancun, May 2000.

- 
- [22] F. GREVE, M. HURFIN, R. MACEDO, M. RAYNAL, « Time and Message Efficient S-Based Consensus », in : *Proc. of the 19th ACM SIGACT-SIGOPS Int. Symposium on Principles of Distributed Computing (PODC'00)*, ACM, p. 332–332, Portland, OR, Jul 2000. Brief announcement.
- [23] J. HELARY, M. HURFIN, A. MOSTEFAOUI, M. RAYNAL, F. TRONEL, « Computing Global Functions in Asynchronous Distributed Systems with Process Crashes », in : *Proc. of the 20th IEEE Int. Conf. on Distributed Computing Systems (ICDCS-00)*, IEEE, p. 584–591, Taipei, Apr 2000. Best paper award.
- [24] J. HELARY, G. MELIDEO, M. RAYNAL, « Tracking Causality in Distributed Systems: A Suite of Efficient Protocols », in : *Proc. of the 7th Int. Colloquium on Structural Information and Communication Complexity (SIROCCO-00)*, p. 181–195, L'Aquila, Jun 2000.
- [25] J. HELARY, G. MELIDEO, « Minimal Size of Piggybacked Information for Tracking Causality: a Graph-Based Characterization », in : *Proc. of the WG'2000 (Workshop on Graphs in Computer Science)*, Konstanz, Jun 2000.
- [26] A. MOSTEFAOUI, M. RAYNAL, F. TRONEL, « The Best of Both Worlds: a Hybrid Approach to Solve Consensus », in : *Proc. of the Int. Conference on Dependable Systems and Networks (FTCS/DCCA)*, IEEE, p. 513–522, New York, NY, Jun 2000.
- [27] A. MOSTEFAOUI, M. RAYNAL, « Consensus Based on Failure Detectors with Perpetual Accuracy Property », in : *Proc of the Int. Parallel and Distributed Processing Symposium (IPPS/SPDP)*, IEEE, p. 514–519, Cancun, May 2000.
- [28] A. MOSTEFAOUI, M. RAYNAL, « k-Set Agreement with Limited Accuracy Failure Detectors », in : *Proc. of the 19th ACM SIGACT-SIGOPS Int. Symposium on Principles of Distributed Computing (PODC-00)*, ACM, p. 143–152, Portland, OR, Jul 2000.
- [29] A. MOSTEFAOUI, M. RAYNAL, « Low-Cost Consensus-Based Atomic Broadcast », in : *Proc of the IEEE Pacific Rim Int. Symposium on Dependable Computing (PRDC'2000)*, IEEE Computer Society Press, Los Angeles, CA, Dec 2000.
- [30] M. RAYNAL, « Logical Instantaneity and Causal Order: Two first class communication modes », in : *Proc. of the European Conference on Parallel Computing (EUROPAR'00)*, A. B. et al. (éditeur), LNCS, 1900, Springer Verlag, p. 35–42, Munich, Aug 2000. Invited paper.
- [31] L. RODRIGUES, R. BALDONI, E. ANCEAUME, M. RAYNAL, « Deadline-constrained causal order », in : *Proc. of the 3rd IEEE International Symposium on Object-Oriented Real-Time Distributed Computing (ISORC 2000)*, IEEE, p. 234–241, Newport, CA, Mar 2000.
- [32] L. RODRIGUES, M. RAYNAL, « Atomic Broadcast in Asynchronous Crash/Recovery Distributed Systems », in : *Proc. of the 20th IEEE Int. Conf. on Distributed Computing Systems (ICDCS-00)*, IEEE, p. 288–297, Taipei, Apr 2000.
- [33] L. RODRIGUES, M. RAYNAL, « Quorum-Based Replication in Asynchronous Crash-Recovery Distributed Systems », in : *Proc. of the European Conference on Parallel Computing (EUROPAR'00)*, A. B. et al. (éditeur), LNCS, 1900, Springer Verlag, p. 605–608, Munich, Aug 2000.

**Rapports de recherche et publications internes**

- [34] E. ANCEAUME, J. HELARY, M. RAYNAL, «Tracking Immediate Predecessors in Distributed Computations», *Research Report n° 1344*, IRISA, Jul 2000.
- [35] F. BRASILEIRO, F. GREVE, M. HURFIN, J. L. NARZUL, F. TRONEL, «Eva: an Event-Based Framework for Developing Specialised Communication Protocols», *Research Report n° 1345*, IRISA, Jul 2000.
- [36] F. BRASILEIRO, F. GREVE, A. MOSTEFAOUI, M. RAYNAL, «Consensus in One Communication Step is Possible», *Research Report n° 1321*, IRISA, Apr 2000.
- [37] P. EZHILCHELVAN, J. HELARY, M. RAYNAL, «Building TMR-Based Reliable Servers Despite Bounded Input Lifetimes», *Research Report n° 1293*, IRISA, Jan 2000.
- [38] P. EZHILCHELVAN, A. MOSTEFAOUI, M. RAYNAL, «Randomized Multivalued Consensus», *Research Report n° 1320*, IRISA, Apr 2000.
- [39] U. FRITZKE, P. INGELS, «Système transactionnel pour données partiellement dupliquées, fondé sur la communication de groupes», *Research Report n° 1322*, IRISA, Apr 2000.
- [40] F. GREVE, M. HURFIN, M. RAYNAL, F. TRONEL, «Primary Component Asynchronous Group Membership as an Instance of a Generic Agreement Framework», *Research Report n° 1292*, IRISA, Jan 2000.
- [41] M. HURFIN, R. MACEDO, A. MOSTEFAOUI, M. RAYNAL, «A Sliding Round Window Diamond-S-Based Consensus Protocol», *Research Report n° 1304*, IRISA, Feb 2000.
- [42] M. HURFIN, A. MOSTEFAOUI, M. RAYNAL, «A Comprehensive Approach for Failure Detector-Based Consensus Protocols», *Research Report n° 1303*, IRISA, Feb 2000.
- [43] A. MOSTEFAOUI, S. RAJSBAUM, M. RAYNAL, «Conditions on Input Vectors for Consensus Solvability in Asynchronous Distributed Systems», *Research Report n° 1355*, IRISA, Oct 2000.
- [44] A. MOSTEFAOUI, M. RAYNAL, «A Condition for k-Set Agreement in Asynchronous Distributed Systems», *Research Report n° 1339*, IRISA, Jul 2000.
- [45] A. MOSTEFAOUI, M. RAYNAL, «Randomized k-Set Agreement», *Research Report n° 1340*, IRISA, Jul 2000.
- [46] M. RAYNAL, «Quiescent Uniform Reliable Broadcast as an Introductory Survey to Failure Detector Oracles», *Research Report n° 1356*, IRISA, Oct 2000.