

Une CNS pour l'acheminement de messages instantanément stabilisant

Alain Cournier¹, Swan Dubois² et Vincent Villain¹

¹Laboratoire MIS, Université de Picardie Jules Verne, 33 rue Saint Leu, 80039 Amiens Cedex 1 (France)
{alain.cournier;vincent.villain}@u-picardie.fr

²LIP6 - UMR 7606/INRIA Rocquencourt, Équipe-projet REGAL, Université Pierre et Marie Curie - Paris 6
104 Avenue du Président Kennedy, 75016 Paris (France), swan.dubois@lip6.fr

Un algorithme *instantanément stabilisant* assure qu'il se comporte toujours conformément à ses spécifications en partant d'une configuration initiale quelconque. Dans cet article, nous nous intéressons au problème de l'acheminement de messages dans un réseau doté d'une commutation de messages. Nous devons gérer les ressources du réseau de manière à pouvoir délivrer des messages à tout processeur du réseau. Dans ce but, nous utilisons l'information fournie par un algorithme de routage. Mais, en raison du contexte de la stabilisation, cette information peut être initialement incorrecte. C'est pourquoi l'existence d'algorithmes instantanément stabilisants pour cette tâche (démontrée dans [CDV09]) implique que nous pouvons demander au système de commencer à acheminer des messages même si les tables de routage sont initialement corrompues. Dans cet article, nous généralisons le résultat précédent en donnant une condition nécessaire et suffisante pour résoudre ce problème de manière instantanément stabilisante.

Keywords: Stabilisation instantanée, acheminement de messages, routage sans interblocage

1 Introduction

De nombreux concepts de tolérance aux pannes ont été introduits en systèmes distribués. Par exemple, l'auto-stabilisation ([Dij74]) assure que le système, indépendamment de son état initial, retrouve un comportement répondant à ses spécifications *en un temps fini* sans intervention externe. L'état initial quelconque peut permettre de modéliser l'effet de fautes transitoires sur le réseau. Un autre concept, la stabilisation instantanée ([BDPV07]), garantit que le système, indépendamment de son état initial, a *toujours* un comportement répondant à ses spécifications. Dans un système distribué, il est classiquement supposé que tout processeur ne peut communiquer directement qu'avec ses voisins. Pourtant, ce processeur peut avoir besoin de communiquer avec tout processeur du réseau. Dans ce but, il y a en réalité deux problèmes à résoudre : la détermination du chemin à suivre par les messages (problème du routage) et la gestion des ressources du réseau réservées au transport des messages (problème de l'acheminement). Ces ressources sont des espaces mémoire (appelés buffers) utilisés pour stocker temporairement les messages durant leur acheminement. Ces deux problèmes ont été largement étudiés (cf. [Tel00, CDV09] pour les références). Nous disposons de nombreuses solutions stabilisantes pour le premier problème. C'est pourquoi dans cet article, nous nous intéressons aux solutions stabilisantes pour le deuxième. [CDV09] démontre qu'il existe des algorithmes d'acheminement instantanément stabilisants (à condition qu'un algorithme de calcul de table de routage auto-stabilisant s'exécute simultanément). Cela implique que nous pouvons demander au système de commencer à acheminer des messages même si les tables de routage sont initialement corrompues. Dans la suite, nous appellerons message valide tout message qui a été généré par un processeur (par conséquent, un message invalide est un message présent dans la situation initiale). Nous pouvons alors spécifier le problème de l'acheminement de la façon suivante : (1) Tout message peut être émis en un temps fini et (2) Tout message valide sera délivré à son destinataire en un temps fini.

L'objectif de cet article est de présenter une condition nécessaire et suffisante pour obtenir un algorithme d'acheminement instantanément stabilisant. Il s'inspire d'un résultat similaire obtenu dans un environnement sans fautes ([MS78, TS81]). La suite de l'article est structurée comme suit : dans un premier temps,

nous présentons le résultat obtenu dans [MS78, TS81] pour un environnement sans fautes (section 2) puis nous donnons notre contribution dans la section 3.

2 Théorème dans un environnement sans fautes

Dans cette section, nous nous plaçons dans un système distribué qui ne peut pas subir de fautes. La configuration initiale est définie : les tables de routage sont correctes et il n'y a aucun message invalide dans les buffers. Dans cet article, nous nous plaçons dans un réseau à commutation de message (cf. [Tel00]). Chaque processeur dispose de $b \in \mathbb{N}$ buffers de taille suffisante pour contenir tout message. La méthode de commutation est composée de trois types de mouvements :

- **Génération** : c'est la "création" d'un nouveau message. Nous supposons que celle-ci est autorisée dès qu'un buffer du processeur émetteur est libre.

- **Transmission** : c'est la copie d'un message dans un buffer du processeur suivant sur le chemin calculé par l'algorithme de routage. Nous supposons que ce mouvement est autorisé dès qu'un buffer sur le processeur en question est libre. En conséquence de ce mouvement, le buffer initial se libère en un temps fini.

- **Consommation** : c'est le mouvement qui libère un buffer occupé par un message à destination du processeur sur lequel est situé ce buffer. Le message est alors délivré à son destinataire. Nous supposons que ce mouvement est toujours autorisé.

Cependant, si nous n'effectuons aucun contrôle supplémentaire sur les mouvements de message, le réseau peut atteindre des situations inacceptables comme des interblocages. Il est alors impossible d'assurer les spécifications du problème. C'est pourquoi, il est nécessaire de définir un algorithme (appelé contrôleur) qui autorise ou interdit dynamiquement (en fonction de l'occupation courante des buffers) certains mouvements. Si la réponse à la question "Est-ce qu'un contrôleur C empêche le réseau d'atteindre un interblocage quelle que soit l'exécution issue de la configuration initiale?" est affirmative, alors C est dit sans interblocage. En assurant l'absence de famine et de perte de messages, nous pouvons constater qu'un tel algorithme répond au problème de l'acheminement de messages. Nous allons présenter une classe de contrôleurs sans interblocage basée sur le concept de *graphe de buffers*, introduit par [MS78]. Plus précisément, la structure choisie est un DAG pour la raison suivante : un interblocage provient du fait qu'il existe un circuit d'attente de libération de buffer. L'idée de base est alors de définir un DAG sur l'ensemble des buffers du réseau de manière à ce que les messages suivent les chemins de ce DAG. Ainsi, il ne peut pas se former de circuit d'attente de libération de buffer. Posons la notation suivante : $G = (V, E)$ est le graphe modélisant le réseau (V est l'ensemble des processeurs et E l'ensemble des liens de communications).

Définition 1 (Graphe de buffers) *Un graphe de buffers $BG = (\mathcal{B}, BE)$ sur un graphe G muni d'un ensemble \mathcal{B} de buffers et d'un ensemble \mathcal{P} des plus courts chemins (induit par l'algorithme de routage) est défini de la manière suivante : (1) BG est un graphe orienté, (2) pour tout chemin $p \in \mathcal{P}$, il existe un chemin dans BG dont la contraction est p^\dagger , (3) pour chaque noeud u de G et pour chaque message m possible, il existe un buffer adéquat[‡] de BG noté $fb(m, u)$ sur u , (4) pour chaque buffer b de \mathcal{B} situé sur un processeur u et pour chaque message m (non destiné à u) possible, il existe un unique buffer adéquat, situé sur u ou sur un de ses voisins, noté $nb(m, b)$ et (5) BE est l'ensemble des arcs $(b, nb(m, b))$ pour tout buffer $b \in \mathcal{B}$ et pour tout message m (non destiné au processeur sur lequel est b) possible.*

La notation $fb(m, u)$, signifiant "first buffer", représente le buffer dans lequel est placé le message m généré par le processeur u . De même, la notation $nb(m, b)$, signifiant "next buffer", représente le buffer dans lequel sera transmis le message m qui occupe le buffer b . Une fois un tel graphe de buffer défini, il est possible de lui associer un contrôleur selon la définition suivante :

Définition 2 (Contrôleur associé à un graphe de buffers) *Etant donné un graphe de buffers $BG = (\mathcal{B}, BE)$ sur un graphe G muni d'un ensemble \mathcal{B} de buffers et d'un ensemble \mathcal{P} des plus courts chemins (induit par l'algorithme de routage), nous définissons le contrôleur C_{BG} de la manière suivante :*

- La génération d'un message m sur un noeud u est autorisé si et seulement si $fb(m, u)$ est libre. m est alors placé dans ce buffer.

[†] un chemin c dans BG est une suite de buffers $c = b_1..b_t$. Notons alors p_i le processeur sur lequel se situe b_i pour tout $i \in \{1, \dots, t\}$. La contraction de c est alors $p_1..p_t$ en omettant les éventuelles répétitions.

[‡] un buffer b est adéquat pour m s'il existe un chemin de b à un buffer du destinataire de m sur BG dont la contraction est dans \mathcal{P} .

- La transmission d'un message m contenu dans le buffer b est autorisée si et seulement si $nb(m, b)$ est libre, m est alors placé dans ce buffer et b est libéré.

Nous sommes à présent en mesure de donner le résultat fondamental suivant ([MS78, TS81]) :

Théorème 1 Soit un graphe de buffers $BG = (\mathcal{B}, BE)$ sur un graphe G muni d'un ensemble \mathcal{B} de buffers et d'un ensemble \mathcal{P} des plus courts chemins. Soit le contrôleur C_{BG} associé à BG conformément à la définition 2. Le contrôleur C_{BG} est sans interblocage si et seulement si BG est un DAG.

Le lecteur pourra trouver des exemples de tels contrôleurs dans [Tel00].

3 Contribution

Dans cette section, nous nous plaçons dans un système distribué sujet à des fautes transitoires. La configuration initiale est donc quelconque : les tables de routage sont incorrectes et il y a des messages invalides dans les buffers. Nous reprenons les notations et définitions introduites dans la section 2. Soit C un contrôleur répondant à la définition 2 pour un graphe de buffers BG . Nous supposons que, pour tout message m occupant un buffer b , le buffer $nb(m, b)$ est calculé au moment de la transmission du message en fonction des tables de routage à cet instant. Nous supposons exister un algorithme \mathcal{A} de calcul de table de routage auto-stabilisant et silencieux. Cet algorithme s'exécute de manière simultanée mais prioritaire par rapport à C (noter que cet algorithme ne nécessite que des communications entre voisins, il n'est donc pas utile qu'il utilise les buffers de communication distante gérés par C). Nous posons les notations suivantes : FB et NB représentent respectivement l'ensemble des buffers $fb(m, u)$ pour tout processeur u et pour tout message possible m associés à BG et l'ensemble des buffers $nb(m, b)$ pour tout buffer b et pour tout message m (non destiné au processeur sur lequel se situe b) associés à BG . Nous donnons à présent notre théorème :

Théorème 2 (CNS pour obtenir un acheminement de messages instantanément stabilisant) C répond au problème de l'acheminement de messages de manière instantanément stabilisante à condition que \mathcal{A} s'exécute de manière simultanée si et seulement si :

- (1) BG est un DAG une fois \mathcal{A} stabilisé.
- (2) C est équitable pour la génération de messages et pour la transmission de messages.
- (3) Pour tout message valide m non délivré durant la stabilisation de \mathcal{A} et quelle que soit l'exécution de C , il existe une copie de m située dans un buffer adéquat pour m une fois \mathcal{A} stabilisé.
- (4) Tous les messages occupant des buffers qui ne leur sont pas adéquats une fois \mathcal{A} stabilisé libèrent tous les buffers de FB et de NB en un temps fini quelle que soit l'exécution de C .
- (5) Pour tout message valide m , C maintient au moins une copie de m tant que celui-ci n'est pas délivré.

Le lecteur peut trouver un algorithme répondant à ces propriétés dans [CDV09]. Nous allons à présent donner les idées de la preuve de ce résultat.

Preuve de la nécessité de la condition Nous souhaitons montrer que si C répond au problème de l'acheminement de messages de manière instantanément stabilisante (à condition que \mathcal{A} s'exécute de manière simultanée) alors C vérifie les cinq propriétés du théorème 2. Pour cela, nous allons raisonner par contraposée. Supposons donc que C ne vérifie pas une des propriétés du théorème 2.

1) Si BG possède un circuit une fois \mathcal{A} stabilisé, nous pouvons construire un circuit d'attente sur un ensemble de messages et donc un interblocage.

2) Si C n'est pas équitable pour la génération ou la transmission, cela signifie qu'il peut mettre un message en famine et donc l'empêcher d'être délivré.

3) S'il existe un message valide m , non délivré pendant la stabilisation de \mathcal{A} , tel qu'aucune copie n'occupe un buffer adéquat pour m une fois \mathcal{A} stabilisé, cela signifie que ce message ne pourra pas être délivré.

4) S'il existe un message ne libérant jamais un buffer de FB ou de NB qui ne lui est pas adéquat une fois \mathcal{A} stabilisé, cela signifie qu'il existe un message qui peut être mis en famine pour sa génération (cas de FB) ou une de ses transmissions (cas de NB). Dans ce cas, ce message ne pourra pas être délivré.

5) S'il existe un message valide m tel que C ne maintienne pas une copie tant qu'il n'est pas délivré, cela signifie que m a été perdu. Ce message ne pourra pas être délivré.

Dans tous les cas, nous constatons que C ne répond pas au problème de l'acheminement de messages de manière instantanément stabilisante, ce qui était le résultat à prouver.

Preuve de la suffisance de la condition Nous souhaitons montrer que si un contrôleur C vérifie les cinq propriétés du théorème 2 alors il répond au problème de l'acheminement de messages de manière instantanément stabilisante (à condition que \mathcal{A} s'exécute de manière simultanée).

Lemme 1 *Si les tables de routage sont correctes dans la configuration initiale et si C vérifie les cinq propriétés du théorème 2 alors il répond au problème de l'acheminement de messages de manière instantanément stabilisante.*

Pour prouver ce résultat, il faut reprendre la preuve du théorème 1 (le lecteur pourra se reporter à la preuve du théorème 5.7 dans [Tel00]), ce qui nécessite la propriété 1 du théorème 2 (en constatant que la propriété du théorème 2.4 permet d'assurer que les messages invalides ne perturberont pas la preuve). Cela permet de prouver que C est sans interblocage dans le cas considéré. Les propriétés 2 et 5 du théorème 2 permettent alors de déduire le résultat.

Lemme 2 *Si les tables de routage sont quelconques dans la configuration initiale et si C vérifie les cinq propriétés du théorème 2 alors il répond au problème de l'acheminement de messages de manière auto-stabilisante à condition que \mathcal{A} s'exécute simultanément.*

Par hypothèse, \mathcal{A} est un algorithme de calcul de table de routage auto-stabilisant et silencieux. Cela signifie que les tables de routage seront construites et correctes en un temps fini (grâce au fait que \mathcal{A} soit prioritaire sur C). Or, C s'exécute de manière simultanée, l'occupation des buffers va donc être modifiée (acceptation de nouveaux messages et acheminement de messages). Nous pouvons appliquer le lemme 1 une fois \mathcal{A} stabilisé. Nous savons donc que C répond au problème de l'acheminement de messages de manière instantanément stabilisante lorsque les tables de routage sont correctes et constantes (ce qui est assuré par le fait que \mathcal{A} est silencieux). C répondra donc à la spécification du problème à partir de cette configuration. C est donc auto-stabilisant.

Lemme 3 *Si les tables de routage sont quelconques dans la configuration initiale et si C vérifie les cinq propriétés du théorème 2 alors il ne détruit aucun message valide sans le délivrer.*

Il s'agit d'une conséquence directe des propriétés 3 et 5.

Les lemmes 2 et 3 nous permettent de déduire trivialement que si C vérifie les cinq propriétés du théorème 2, alors il répond au problème de l'acheminement de messages de manière instantanément stabilisante à condition que \mathcal{A} s'exécute simultanément.

4 Conclusion

Dans cet article, nous avons donné un ensemble de propriétés nécessaires et suffisantes pour obtenir un acheminement de messages instantanément stabilisant. En ce sens, notre résultat généralise celui obtenu dans [MS78, TS81]. Il serait possible d'appliquer notre résultat pour démontrer plus simplement l'algorithme de [CDV09].

Références

- [BDPV07] Alain Bui, Ajoy Kumar Datta, Franck Petit, and Vincent Villain. Snap-stabilization and pif in tree networks. *Distributed Computing*, 20(1) :3–19, 2007.
- [CDV09] Alain Cournier, Swan Dubois, and Vincent Villain. A snap-stabilizing point-to-point communication protocol in message-switched networks. In *IPDPS (to appear, preprint available at <http://pagesperso-systeme.lip6.fr/Swan.Dubois/pdf/CDV09a.pdf>)*, 2009.
- [Dij74] Edsger W. Dijkstra. Self-stabilizing systems in spite of distributed control. *Commun. ACM*, 17(11) :643–644, 1974.
- [MS78] Philip M. Merlin and Paul J. Schweitzer. Deadlock avoidance in store-and-forward networks. In *Jerusalem Conference on Information Technology*, pages 577–581, 1978.
- [Tel00] Gerard Tel. *Introduction to Distributed Algorithms*. Cambridge University Press, Cambridge, UK, 2nd edition, 2000.
- [TS81] Sam Toueg and Kenneth Steiglitz. Some complexity results in the design of deadlock-free packet switching networks. *SIAM J. Comput.*, 10(4) :702–712, 1981.