

## 任意形状の非同期ネットワークにおける耐故障リーダー選挙について

松本 哲也 若林 真一 小出 哲士 吉田 典可

広島大学工学部  
東広島市鏡山1丁目4-1

あらまし 分散アルゴリズムの基本的問題の1つにリーダー選挙問題がある。本研究では、リンク故障を含むと仮定した任意形状の完全非同期ネットワークにおけるリーダー選挙問題について考察を行う。まず、各プロセスが固有の識別子を持ち、ネットワークサイズと故障リンク数の上限値を初期情報としてもつとしたとき、リーダー選挙問題に対するメッセージ数 $O(m+n\log^2 n)$ 、メッセージサイズ $O(\log n)$ ビットの終了検知を伴うアルゴリズムを提案する。ここで $n$ はネットワーク中のプロセス数、 $m$ はリンク数である。次に、各プロセスがネットワークサイズを初期情報としてもつとしても固有の識別子をもたなければリーダー選挙問題は解けず、たとえ各プロセスが独立に確率的選択を行うランダム化アルゴリズムを用いても解けないことを示す。

和文キーワード 分散アルゴリズム, リーダ選挙, 非同期システム, リンク故障, 匿名ネットワーク, 終了検知

## On Leader Election for Faulty Asynchronous Networks with Arbitrary Topology

Tetsuya MATSUMOTO Shin'ichi WAKABAYASHI Tetsushi KOIDE Noriyoshi YOSHIDA

Faculty of Engineering, Hiroshima University  
4-1, Kagamiyama 1 chome, Higashi-Hiroshima, 724 JAPAN

**Abstract** Leader election is one of the basic problems on distributed networks. In this paper, leader election for fully asynchronous distributed networks with arbitrary topology with link faults is considered. First, when every process has a unique identifier and knows the size of network and the upper bound of number of link faults as initial information, we propose an algorithm for leader election with termination detection that uses  $O(m+n\log^2 n)$  messages, each of size  $O(\log n)$  bits, where  $n$  and  $m$  are the number of processes and links, respectively, in the network. Next, when every process knows the size of network but doesn't have unique identifiers, we show that even if a randomized algorithm that can execute random selection locally is used, leader election cannot be solved.

英文 key words distributed algorithm, leader election, asynchronous system, link fault, anonymous network, termination detection

## 1. まえがき

分散アルゴリズムの基本的問題の1つにリーダー選挙問題 (LEP) がある<sup>[10][11]</sup>。これは、トークンシステムにおいてトークン消失後にトークンを再生成させるプロセスを選択する場合など、分散システム中で特別な役割をもたせるプロセスをただ1つ選択することを目的とした問題であり、様々な応用がある。一方、分散システムでは耐故障性の考慮が重要な課題である。本稿では特に、故障を含む任意形状のネットワークにおけるLEPについて考察する。従来、この問題は各プロセスが固有の識別子をもつと仮定して解かれており、故障を含み、かつ固有の識別子をもたないネットワークにおけるLEPについての研究はほとんど行われていない。本稿では固有の識別子をもつ場合とまたない場合のLEPについて考察を行う。

## 2. 準備

本稿で考察する分散ネットワークは通信リンクを介して通信するプロセスの集合であると仮定する。各プロセスは独自に同じアルゴリズムを実行し、共有変数を持たず、プロセス間の通信は2プロセス間にリンクが存在するときのみ直接メッセージを送ることができる。ネットワーク全体は完全非同期である。ネットワークは $n$ 個のプロセスと、 $m$ 本の双方向通信リンクの集合からなり、プロセスをノード、リンクを枝とする任意形状のグラフで表される。またアルゴリズム開始時、任意のプロセス部分集合が自発的に起動する始動プロセスに成り得るとする。各プロセスは初期情報として、実行するアルゴリズム、各自のポートをもち、ポートは区別できるとする。

本稿で対象とする故障はリンク故障である。故障リンクでは送信が行われても受信側プロセスにメッセージを伝えないものとする。すべての故障はアルゴリズム開始時に既に起こっているとし、プロセス故障は考えず、ネットワークは非故障リンクのみで連結とする。

本稿においてLEPを解くとは、

(A) 1. ただ1つの特別なプロセスがリーダーとして選ばれ、

2. 他のすべてのプロセスはリーダーを知ることであるとする。ただし、LEPを解くアルゴリズムはメッセージドリブンとする。

識別子と故障に関する条件を変えた種々のグラフとそれらのグラフに対するLEPに関する過去の研究結果を表1にまとめる。ネットワークが固有の識別子をもたない場合、故障やネットワークサイズの知識にかかわらずLEPを決定的に解くことができない場合が存在することが知られている<sup>[2]</sup>。従って本稿では、アル

ゴリズムが終了しても問題が正しく解けなかった場合でも、各プロセスが解けなかったことを検知できるなら解けたとみなす。これは、もし解けなかったことが検知できるならランダム化アルゴリズムを繰り返し実行すればいずれ必ず条件(A)でLEPは解けるので許容できるといえる。結局、本稿においてLEPを解くとは条件(A)を満たすか、次の条件(B)を満たせばよいとする。条件(B)で解ける場合は表1中では三角形で示している。

(B) アルゴリズム終了時に(A)の条件を満たさない場合、各プロセスはそのことを検知する。

表1 種々のネットワークに対するLEPについての過去の研究結果

	識別子○		識別子×	
	n未知	n既知	n未知	n既知
故障無	G ○	G <sub>n</sub> ○	G <sub>n</sub> × <sup>[2]</sup>	G <sub>m</sub> △ <sup>[5-8]</sup>
故障有	G <sub>f</sub> × <sup>[4]</sup>	G <sub>n</sub> ○ <sup>[1][3]</sup>	G <sub>f</sub> ?	G <sub>m</sub> ?

識別子○：固有の識別子をもつネットワーク

識別子×：固有の識別子をもたないネットワーク

故障無：故障を含まないと仮定したネットワーク

故障有：故障を含むと仮定したネットワーク

n既知：各プロセスはnを初期情報としてもつ

n未知：各プロセスはnを初期情報としてもたない

○：必ず条件(A)を満たす

△：条件(A)または(B)を満たす

×：条件(A)も(B)も満たさない

ネットワークのクラスGの添え字のfは故障を含むこと(faulty), aは固有の識別子をもたないこと(anonymous), nはネットワークサイズ既知であること(number of nodes)を表す。例えば、固有の識別子をもたず、故障を含み、各プロセスがネットワークサイズnを既知であるようなネットワークのクラスはG<sub>m</sub>で表す。

本稿ではG<sub>n</sub>とG<sub>m</sub>について考察を行った。[1]ではG<sub>n</sub>においてメッセージ数 $O(m+n \log^5 n)$ でLEPを解くアルゴリズムが提案されている。G<sub>n</sub>の各プロセスが故障リンク数の上限値 $f_i$ (定数)を初期情報としてもつとしたネットワークのクラスをG<sub>n</sub>'と定義する。3. ではG<sub>n</sub>'に対して条件(A)を必ず満たしてメッセージ数 $O(m+n \log^2 n)$ でLEPを解くアルゴリズムALEPを提案する。また、4. ではG<sub>m</sub>に対しLEPは解けないことを示す。

## 3. 固有の識別子をもち、故障を含むネットワークについて

本節ではネットワークG<sub>f</sub> ( $\in G_{n'}$ ) 上でLEPを解くアルゴリズムALEP(Algorithm for LEP)について示す。

### 3. 1 アルゴリズムALEP

ALEPはスパニング森を成長させる手続きGSTと、スパニング森の連結成分である各木のノード数を数える手続きNCからなり、これらは同時に実行される。GSTによって始動プロセスから根付きスパニング森を構築していく。ところが $G_i$ は故障リンクを含むのでGSTでは木のノード数は不明で終了検知ができない場合が存在する。よってNCにより木のノード数を数え、 $n/2$ 以上のノードを含む木ができたとき、そのルートがリーダーとなる。以下ではGSTとNCを分けて説明する。

#### 3.1.1 スパニング森生成手続きGST

GSTは[1]の終了検知なしの耐故障リーダー選挙手法とほとんど同じものを用いる。GSTで構成するスパニング木のリンクをBranchと呼ぶ。Branchの集合は森を構成する。森の各連結成分をフラグメント、フラグメントのルートのコアと呼ぶ。最初Branchの集合は空で、各始動プロセスがそのフラグメントのコアになる(図1)。

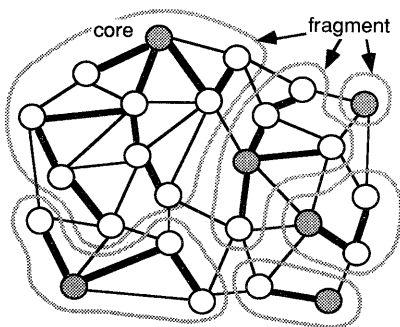


図1 コアとフラグメント

GSTにおいて、各フラグメントの動作はコアによって制御され、隣接する他のフラグメントと結合するためのリンクを1本選択する。選ばれたリンクはBranchになってフラグメント同士が結合し、どちらか一方のフラグメントのコアが拡張されたフラグメントのコアになる。いずれただ1つのフラグメントが残り、そのコアがリーダーになる。

フラグメントは他のフラグメントと結合しながら成長するが、他のフラグメントと識別するためにラベルを用いる。各フラグメントは(レベル, コアの識別子)の順序対であるラベルをもつとし、レベルは非負の整数であるとする。アルゴリズム開始時、各フラグメントは1ノードからなり、レベルは0とする。

#### (A) リンクの種類

ネットワークは非同期のため、故障リンクを遅延の大きなリンクと区別することは不可能であり、アルゴリズムは各リンク上でメッセージが必ず有限時間内で通信できることを仮定して動作することは許されない。それ故、各ノード $v$ は隣接リンクを、メッセージを受け取ったか否かによりOperational, Quietに分類する。

初期時、全リンクはQuietとする。ノードは自発的、もしくはメッセージ受信により起動すると、まずHelloメッセージを隣接リンクの各々に送る。Helloメッセージを受け取ったノードはそのリンクをOperationalと指定する。リンクがQuietであるかぎりは、さらなるメッセージは送らない。

Operationalリンクは以下のように3つに分類される(図2)。

- (i) Branch( $v$ ) : フラグメントの木のリンクでノード $v$ に接するリンク集合  
 $v$ がコアでなければコアの方向にparent( $v$ ), 逆向きにchildren( $v$ )と設定する。
- (ii) Rejected( $v$ ) : 同じフラグメントの他のノードに接続し、Branchではないリンク集合  
 いったんRejectedと分類されれば存在しないものとして扱われる。
- (iii) Linkqueue( $v$ ) : (i), (ii)以外のリンク集合  
 Operationalになるとキューに加えられる。キューの先頭のリンクをtestlink( $v$ )とし、ノードはtestlink( $v$ )をBranch( $v$ )に加えようとする。

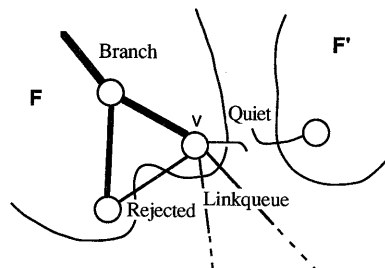


図2 リンクの種類

#### (B) 外向枝の選択

フラグメントのコアがそのラベル(レベル)を変えるときはいつでも、

- (1) 新しいラベルをフラグメント内の各ノードに知らせ、
  - (2) 新しい外向枝選択フェーズを始める、
- ために、Initiateメッセージをフラグメントの全ノードに放送する。ノードはInitiateメッセージを受け取ると

そのラベルを更新し、メッセージを受け取ったリンクをparent(v)とし、同じメッセージを子に放送する。このラベル更新手続きをINIT手続きと呼ぶ。それから、後に述べるテスト動作により、Linkqueue中から自フラグメントより大きいラベルのフラグメントを導く外向枝を探し始める。最初に見つかった枝をOutlink(v)とし、親にPossibleOutlinkメッセージを送る。vはLinkqueue中からどれか1つをOutlink(v)として見つかるか、もしくは子からPossibleOutlinkメッセージを受け取るかして外向枝が見つかったことを知る。vの子もまた同じことをする。いったんOutlink(v)を選ぶと、その後で報告されたものは無視される。

リンクの探索はTestメッセージをtestlink(v)に送ることで行う。ノードvがノードwにTestメッセージを送ったとする(図3(a))。

- (1) vとwが同じラベルをもつとき (図3(b))  
(すなわち同じフラグメントに属する)

wはvにRejectメッセージを返し、vとwはそのリンクをRejectedと指定する。そしてvはLinkqueue(v)中の次のリンクを新しいtestlink(v)とする。

- (2) wのラベルの方が大きいとき (図3(c))  
wはvにPossibleOutlinkメッセージを返す。

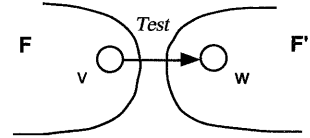
- (3) wのラベルの方が小さいとき (図3(d))  
とりあえず返答しない。いつかwが大きくなったとき返答するか、もしくはwがvにTestメッセージを出すことになる。(ただし、NCのために、wはTestを受け取り、wのラベルの方が小さかった事実をvに伝えることだけは行う。)

vはPossibleOutlinkメッセージをwから受け取ると、もし以前に受け取っていなければ、そのtestlink(v)をOutlink(v)とし、親へ伝え、いずれこのことはコアまで伝わることになる(図3(e))。

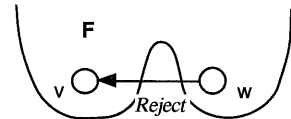
Initiateメッセージに関連して使われるすべてのメッセージはフラグメントのラベルでタグ付けされている。新しいInitiateメッセージを受け取ると現在のリンク選択手続きはアボートされる。すなわち、ノードはOutlinkをnullにリセットし、前のレベルのPossibleOutlinkメッセージは無視する。

### (C) フラグメントの結合

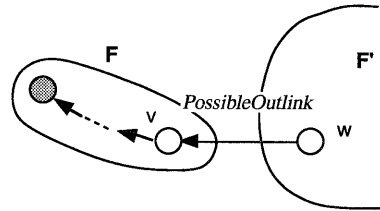
いったんコアがOutlinkを選ぶと、Outlinkに沿ってConnectメッセージが送られ、そのOutlinkに対応する外向枝を見つけたvに伝えられる。vはそのOutlinkをBranchに変え、リンクのもう片方の端点であるwにConnectメッセージを送る。フラグメントの結合には次の2通りがある。



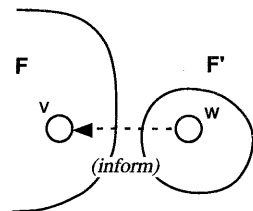
(a) Testメッセージ



(b) label(v)=label(w)のとき



(c) label(v)<label(w)のとき



(d) label(v)>label(w)のとき

図3 リンクのテスト

- (1) level(v)<level(w)のとき

[吸収手続き] (図4(a))

レベルの小さい方のフラグメントに属するノードのみがラベルを更新し、大きい方のラベルを受け継ぐ。具体的には、wはvにInitiateメッセージを送る。すなわちwのラベルがv側のフラグメントに放送されることになる。そしてv側のフラグメントはOutlinkを探すwのフラグメントサーチに参加する。このとき、w側のフラグメントではv側のフラグメントが追加されたことを知らせるメッセージは使われない。

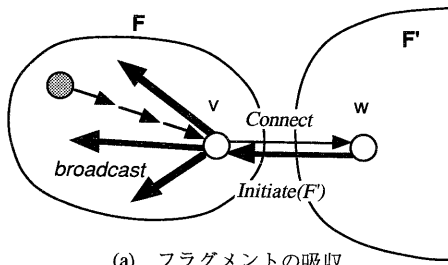
- (2) level(v)=level(w)=jのとき

[合併手続き] (図4(b))

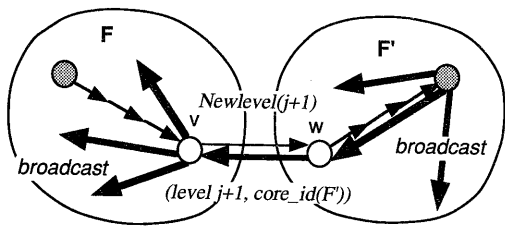
両方のフラグメントに属するノードがレベルをインクリメントし、どちらかのコアの識別子をもつようラベル更新する。具体的には、wはNewlevel(j+1)メッセージを親に送り、コアまで伝える。もし途中

のノードで他のNewlevelメッセージやInitiateメッセージが先に受け取られていれば、このNewlevelメッセージはアボートされ、これ以上何もしない。そうでなければメッセージはコアに到着する。もしコアのレベルが既に上がって居れば、もしくは既にConnectメッセージを送って居ればこのNewlevelメッセージは無視される。そうでなければレベルをj+1にし、新しいフラグメントラベルでInitiateメッセージを全ノードに放送する。これはいずれwに到着し、いずれvのフラグメントのノードすべてに到達する。

Newlevelメッセージがアボートされるときは、wのフラグメントのレベルが他の理由で増加していることになり、それ故にいずれInitiateメッセージが受け取られることになる。



(a) フラグメントの吸収



(b) フラグメントの合併  
図4 フラグメントの結合

レベルjのフラグメントは少なくとも $2^j$ のノードをもつので、可能な最大レベルは $\log n$ となる。ノードのフラグメントレベルは単調増加であり、あるノードが同じレベルの間はそのフラグメントのコアの識別子は変わらないので、各ノードは高々 $\log n$ 個のフラグメントに属することになる。

### 3.1.2 ノード数計算手続きNC

フラグメントのノード数が $n/2$ を超えるとそのコアはリーダとなり、そのことを放送してアルゴリズムは終了する。しかし、GSTの動作のみではノード数が不明で終了検知ができない。よってNCではフラグメント内のノード数を数えることを目的とする。レベルjのフラグメントには $2^j$ 以上のノードが含まれている。NC

Cではこのフラグメントのノード数が $2^{j+1}$ を超えたとき必ずレベルj+1のラベルに更新する。こうすると過半数のノードを含むフラグメントは必ず $j \geq \log n - 1$ なので、ここで1つでもノードが加わるとコアはリーダとなり、アルゴリズムが終了する。以下で具体的な動作について説明する。

あるフラグメントのレベルが上がると、ノードのラベルが更新されると、GSTの再開と同時にNCが開始される。各フラグメントに対し、 $2f_i + 1$ 個のトークンがコアから深さ優先順にノード数を数えながらBranchから構成される木を探索する(図5)。このメッセージにはフラグメントの前レベルまでの正確なノード数が含まれている。リンクには、Branch, Rejected, Linkqueue, Quietの高々4種類が存在する。あるノードにおいてトークンが次に探索するリンクの種類により、NCは以下の動作を行う(図6)。

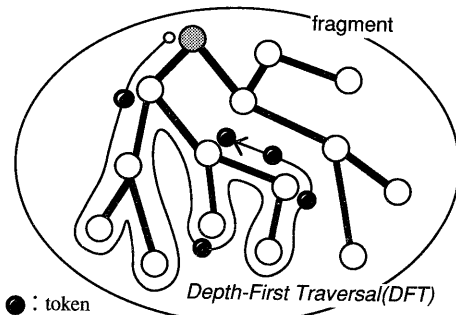


図5 トークンによる深さ優先探索

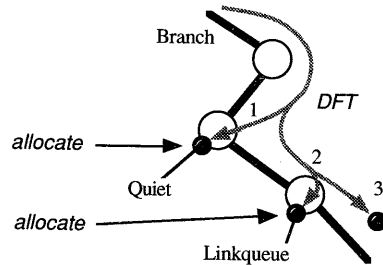


図6 トークンの探索と割り当て

#### Case 1. Branch

探索順どりにトークンを進める。

#### Case 2. Rejected及び既にトークンの割り当てられたリンク

そのリンクを無視し、そのノードの次のリンクの探索に移る。

#### Case 3. Quiet

まだトークンが割り当てられていなく

ればトークンを割り当てる。既に割り当てられていれば、トークンは次のリンク探索に向かう。この後、このリンクがOperationalに変わればCase 4.の動作を行う。

**Case 4. Linkqueue**

まず、そのリンクがRejected, 外向枝のいずれかに分類されるまで待ち, RejectedになればCase 2.の動作に従う。外向枝と判断された場合は1トークンを割り当てる。後の動作は以下で詳しく説明する。

もしトークンがコアに戻って来たときのノード数が次レベルの必要値を超えていたら, GSTのラベル更新手続きINITを起動する。

Case 4.で探索リンクが外向枝と判断された場合について考える。あるフラグメントFのノードXにおいて1トークンがリンクLに割り当てられ, Lに接するもう一方のフラグメントをF'とする。次の2つの場合を考える。

**(1)level(F)≤level(F')のとき**

GSTの動作より, レベルが等しければ合併手続き, 異なれば吸収手続きにより, いずれにしてもFのラベル更新が行われる。よってこのレベルのNCは自動的にアポットされる。

**(2)level(F)>level(F')のとき**

Lに接するF'側のノードをYとする。GSTの動作より, XはF'のレベルがFのレベルより小さいことを知る。するとXはLにトークンを割り当て, これ以降Xに到着したトークンはLをRejectedと同様に扱う。F'のレベルはFより小さいので, 今後F'がFに吸収されることがあり, この場合吸収したフラグメントのノード数も数える必要がある。吸収が起こった場合の動作は以下のようなになる。

吸収の際, XはF'のノード数 $n_{F'}$ を知ることができる。よってF'のノードが加わるにより次レベルの必要ノード数を超える場合は, Fのコアに手続きINITを起動させるようメッセージを送る。超えていなくても他のリンクを通して別のフラグメントが吸収されている可能性があり, 全体のノード数を $n_F + n_{F'} \times (2f_1 + 1)$ と見積もる (図7)。

**[1]見積もりが次レベルの必要ノード数を超える場合**

この場合, 放送によるフラグメント内ノード数計算手続き(CNB)を起動する。すなわち次レベルに必要なノード数の残りの $1/(2f_1 + 1)$ が加わればCNBを起動することになる。CNBの動作は以下の通りで, Branchリンクのみを使って行われ

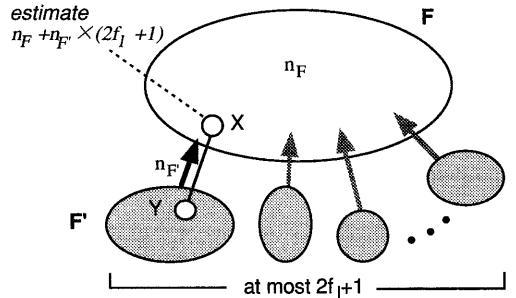


図7 吸収の見積もり

る。まずコアからメッセージCNodesを現在フラグメントに属している全ノードに放送する。そしてリーフノードは親ノードにノード数1を返し, 内部ノードは子ノードすべてから返されたノード数の合計に1を足して親ノードに返す。コアは子のすべてからノード数を受け取ると, それを合計してフラグメントの実際のノード数を知ることができる (図8)。

そして実際にフラグメントのラベル更新に必要なノード数を超えていれば手続きINITをFのコアが起動する。超えていなければ現レベルのまま, 実ノード数をもってNCを再起動する。このとき, F'はFに吸収されており, 新しく成長したFで探索が行われる。

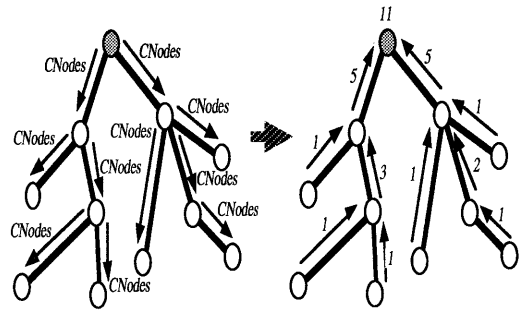


図8 手続きCNBの動作

**[2]見積もりが次レベルの必要ノード数を超えない場合**

YをF'部分の仮のルートとしてNCを再帰的に起動する (図9)。すなわち元F'部分上で $2f_1 + 1$ トークンにより深さ優先順で探索し, XがCNBを起動するために必要なノード数をYが集めなければならないノード数とする。元F'中のノードはその $1/(2f_1 + 1)$ を集めればF'部でCNBを起動し, これを繰り返すことによりYの必要数に達すればXに伝える。

Branch以外の全リンクがRejectedになると部分木の探索は終了する。Xは新しく加わったノード数をトークンに記録して、このトークンはFの次リンク探索へと進む。訪れたノードではノード数を更新していく。探索の途中で次レベルの必要ノード数に達したときも手続きINITを起動する。

また、トークンは深さ優先探索を行うので1パスで同じノードを2回訪れることになるが、2回目は何も行わず、Branchを使って移動のみを行う。

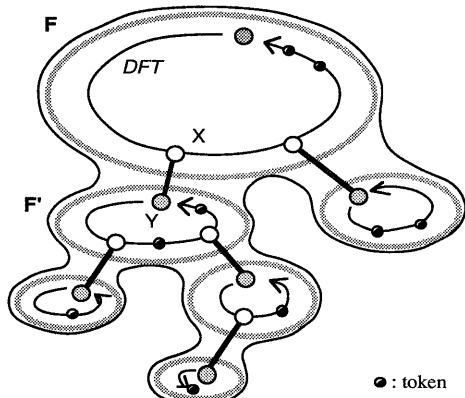


図9 NCの再帰的起動

以上がALEPの動作である。プロセス故障を含むネットワークに対しても、故障したプロセスに接続している全リンクが故障リンクであると考えることにより、ALEPが適用可能である。

### 3. 2 ALEPの正当性の証明

本アルゴリズムはノード数が $n/2$ 以上になったフラグメントのコアがリーダ宣言し、全ての処理を強制的に終了させる。従って証明すべき点は、

1. GSTはいずれ1つのスパンニング木を構成する。
2. NCはフラグメントのノード数をコアにおいて正しく計算する。

の2つである。1. に関しては[9]などから明らかなので省略する。

GSTのラベル更新手続きINITを行う際にルートは自フラグメントのノード数を数えることができるので、2. は以下のように置き換えることができる。

- 2'. フラグメントのノード数が次レベルに必要なノード数に達したときは、いずれ必ず手続きINITを起動する。

以下ではこれを証明する。

あるフラグメントFを考え、その上でNCを実行することを考える。Fのコアをr、ノード数 $n_F$ とする。

Fに属する枝の種類はBranch, Rejected, Linkqueue, Quietである。トークンtはF上をrから深さ優先順に探索する。次に探索するリンクがBranchなら深さ優先順に次にノードへ進む。Rejectedなら無視し、そのノードの次のリンクへ探索を進める。Quietで、まだトークンが割り当てられていないリンクなら、トークンを1つそのリンクに割り当てる。既に割り当てられていればトークンはそのノードの次のリンク探索に進む。従って、もしFの全リンクがBranch, Rejected, Quietのいずれかで、Quietリンクはすべて故障リンクであるとしても、故障リンク数の上限は $f_1$ であり、使用トークン数 $2f_1+1$ だから少なくとも1トークンはいずれ必ずrに帰ってくる。よってrは $n_F$ を知ることができ、もし次レベルに必要なノード数に達していれば手続きINITを起動する。

次に探索リンクがLinkqueueであるときを考える。これはトークンを割り当てられたQuietリンクが後にLinkqueueになった場合も含む。このリンクLにはGSTによりTestメッセージが送られ、Rejectedか外向枝に判断される。Rejectedであれば上記と同様、そのノードの次のリンク探索に進む。外向枝と判断された場合、以下の2つの場合が考えられる。

相手側のフラグメントをF'とする、最初の場合は $label(F) \leq label(F')$ のときである(図3(c))。このときはLをPossibleOutlinkとしてrに伝えようとする。もし伝わればLをOutlinkとしてFはF'に吸収されるか、もしくは合併が起こる(図4(a), 4(b))。すなわちFはラベル更新を行うので、その際に $O(n)$ のメッセージ数でFの正しいノード数を計算でき、このレベルにおけるトークンによる探索はアボートしてよい。またrがLをOutlinkと選ばなくても、別の外向枝をOutlinkとして吸収または合併が起こるので、同様にこのレベルのトークンによる探索はアボートしてよい。

2番目の場合は $level(F) > level(F')$ のときである。Lに接するF側のノードをX、F'側のノードをYとする。GSTの動作より、XはF'のレベルがFのレベルより小さいことを知る。するとLにトークンを割り当て、これ以降Xに到着したトークンはLを無視する。F'のレベルはFより小さいのでF'がFに吸収されることがあり、この場合を詳しく考える(図7)。F'はラベル更新を行うのでその際にF'のノード数 $n_{F'}$ をXは知ることができる。Xは現在の $n_F$ に $n_{F'}$ を加えて、もし次レベルに必要なノード数に達するならば手続きINITを起動するようrに要請する。もし達していなくても他のリンクを通して別のフラグメントが吸収されている可能性があり、全体のノード数を見積もる必要がある。トークンの数は $2f_1+1$ 個なのでそうした吸

取は高々 $2f_i+1$ 箇所では起こりえない。たとえ1トークンによって何回も吸収が行われていてもその合計はCNB起動のための必要数を超えていない。従って、全体のノード数を $n_F+n_{F'} \times (2f_i+1)$ と見積もり、これが次レベルのノード数を超えていればフラグメント内ノード数計算手続きCNBを $r$ に起動させる。そして得られた実際のノード数が次レベルの必要数を超えていれば手続きINITを起動する。見積もりでは超えたが実際には超えていなかった場合、得られたノード数を新しい $n_{F'}$ として新たにNCを起動する。見積もりは $2f_i+1$ 倍、すなわちトークンをもつ全ノードでノード数が増加していると考えて求めているので、必要数の残りの $1/(2f_i+1)$ 以上のノードが新しく加わったところが存在するならば必ずCNBが起動される。

また、見積もりが超えなかった場合でも $Y$ を仮のルートとして再帰的にNCを起動するので、後にCNB起動に必要なノード数が増加すれば $X$ は知ることができ、しきい値に達すれば必ずCNBが起動される。すなわち、 $F$ のノード数が次レベルに必要なノード数に達した場合は必ず手続きINITが起動されることになる。

また、 $F'$ の全リンクがRejectedとなった場合、トークン $t$ は次のリンク探索に進むので全トークンが最大 $f_i$ 本のリンクの両端に割り当てられたとしても少なくとも1トークンは動作状態にあり、NCがデッドロックに陥ることはない。

よって $2'$ が証明され、アルゴリズムALEPは正しくLEPを解くことが示された。

### 3. 3 ALEPのメッセージ複雑度の解析

GSTのメッセージ複雑度は $O(m+n \log n)$ であることは[1]などから明らかなので、NCについてのみ解析を行う。

各レベルで考える。まず、あるフラグメント $F$ においてNCが起動されると、 $2f_i+1$ 個(定数)のトークンが $F$ の深さ優先探索を行うので $(2f_i+1) \times 2 \times O(n)$ のメッセージが使われる。すなわち、トークンによる探索は $O(n)$ メッセージが必要である。

トークンを割り当てられたリンクを通して $F'$ を吸収すれば、 $F'$ においてもトークンによる探索が行われる。吸収した元のフラグメント部分のノード数をそれぞれ $n_i$ とすると、 $\sum n_i < n$ であり、各フラグメント探索は重なりがないので、見積もりによるトークン探索が行われなにかぎりは $O(n)$ メッセージしか必要としない。

従って、ワーストケースは吸収で必要ノード数に達さず、見積もりによりCNBを起動しても最小ノード

数しかノードが加わっていない場合が続くときである。あるフラグメント $F_0$ が $F_1$ を吸収し、その後 $F_1$ のリンクを通して $F_i$ を吸収したとする( $F_i$ は複数存在する)。 $F_i$ について考える。

$F_i$ が $F_j$ を吸収することで、 $F_j$ ではそれぞれトークンが回されている。 $(F_0$ の見積もりのための) $F_1$ の必要ノード数の $1/(2f_i+1)$ のノードが $F_i$ に加われば、 $F_j$ でCNBが起動する。CNB起動ごとに必要残り数は減っていくので、これを $O(\log n)$ 回繰り返せば $F_j$ は必ず必要数を得ることができ、 $F_0$ のCNBを起動させることができる。よって、 $F_j$ は必要ノード数を得るために、見積もりによるCNB起動を高々 $O(\log n)$ 回行う。 $F_0$ のCNBを起動させた後、 $F_i$ は $F_0$ 、 $F_i$ を含む新たな深さ優先探索木の一部となる。

NCの動作は、吸収したフラグメントにおいて再帰的にNCを起動することがあるので、部分フラグメントが見積もりによるCNBの起動を $O(\log n)$ 回行って上のレベルの部分フラグメントの必要見積もりノード数を集めることが、鎖状になったフラグメントの数だけ起こる可能性がある。従って1つのフラグメントにおいて1レベルの間に最も多くのメッセージが使われるのは以下のような場合である。

$(\log n-1)$ レベルになると、あとは $O(n)$ メッセージでリーグ宣言し終了するので、1レベルの間にレベル $(\log n-2)$ のフラグメントがレベル $(\log n-3)$ のフラグメントを吸収し、このレベル $(\log n-3)$ のフラグメントがレベル $(\log n-4)$ のフラグメントを吸収し、これが連続してレベル2, 1, 0と鎖上に連続吸収が行われた場合にメッセージ数最大となる(図10)。このフラグメントの各部分を $F_{(\log n-2)}$ ,  $F_{(\log n-3)}$ , ...,  $F_2$ ,  $F_1$ ,  $F_0$ とする。 $F_{(\log n-2)}$ において最多メッセージが使われるのは、上の議論より $F_{(\log n-2)}$ 上でCNBが $(\log n-2)$ 回起動されるときで、 $F_{(\log n-2)}$ のノード数 $n_{F_{(\log n-2)}}$ とすると $n_{F_{(\log n-2)}} \times (\log n-2)$ のメッセージが $F_{(\log n-2)}$ で使われる。次に、 $F_{(\log n-3)}$ において最多メッセージが使われるのは、 $F_{(\log n-3)}$ 上でCNBが $(\log n-3)$ 回起動されることにより $F_{(\log n-2)}$ の見積もりが必要数を超えて $F_{(\log n-2)}$ のCNBを起動するが、実際には必要数を満たさないときである。いったん $F_{(\log n-2)}$ のCNBを起動すると $F_{(\log n-3)}$ は $F_{(\log n-2)}$ の探索木の一部となるので、 $F_{(\log n-2)}$ と結合した後のCNB起動を考えて、 $F_{(\log n-3)}$ で使われる最多メッセージ数は $n_{F_{(\log n-3)}} \times ((\log n-3) + (\log n-2) - 1)$ となる。同様に、 $F_i$ において使用される最大メッセージ数は、

$$\begin{aligned} & n_{F_i} \times \{i + ((i+1)-1) + ((i+2)-1) + \dots \\ & \quad + ((\log n-3)-1) + ((\log n-2)-1)\} \\ & = n_{F_i} \times \{(\log n-3) + (\log n-4) + (\log n-5) + \dots \\ & \quad + (i+2) + (i+1) + i \times 2\} \end{aligned}$$



となる。部分フラグメントの鎖の長さの最大は $\log n - 1$ である。これが起こるのは $n_{F(\log n - 2)}, n_{F(\log n - 3)}, \dots, n_2, n_1, n_0$ がそれぞれ $n/2, n/4, \dots, 8, 4, 2$ より小さいときである。もしいずれかがこれらの値より大きければ途中でラベル更新が起こるはずである。従って、1レベルにおける最大メッセージ数は、

$$\begin{aligned}
 & n/2 \times (\log n - 2) \\
 & n/2 \times (\log n - 3) \\
 & n/8 \times (\log n - 3) + n/4 \times (\log n - 4) \\
 & n/16 \times (\log n - 3) + n/16 \times (\log n - 4) + n/8 \times (\log n - 5) \\
 & n/32 \times (\log n - 3) + n/32 \times (\log n - 4) + n/32 \times (\log n - 5) + n/16 \times (\log n - 6) \\
 & \quad \vdots \\
 & 8 \times (\log n - 3) + 8 \times (\log n - 4) + 8 \times (\log n - 5) + \\
 & 4 \times (\log n - 3) + 4 \times (\log n - 4) + 4 \times (\log n - 5) + \dots \\
 & + 2 \times (\log n - 3) + 2 \times (\log n - 4) + 2 \times (\log n - 5) + \dots \\
 & \hline
 & < n/2 \times (\log n - 2) + \\
 & n \times (\log n - 3) + n/2 \times (\log n - 4) + n/4 \times (\log n - 5) + n/8 \times (\log n - 6) \\
 & + \dots \\
 & < n \log n \times (1/2 + 1/4 + 1/8 + \dots) \\
 & = 2.5 \times n \log n \\
 & = O(n \log n) \\
 & \text{となる。}
 \end{aligned}$$

また、複数のフラグメントが存在する場合を考える。上の議論でフラグメントの最大レベルを1下げると、各部分フラグメントのノード数は半分以下、レベル数も1ずつ下がるので最大メッセージ数は先の半分以下である。こうしたフラグメントの数は高々2倍にしかならないので、いずれの場合も1レベルで使われるメッセージは $O(n \log n)$ の上限を超えることはない。

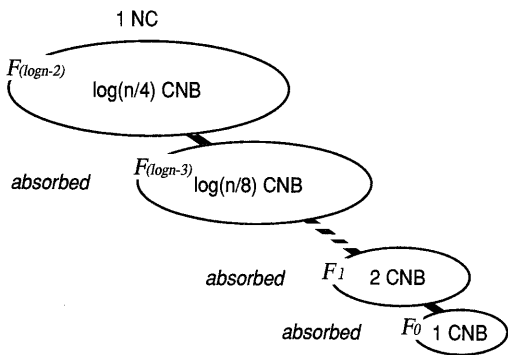


図10 最大メッセージ数を要するフラグメント吸収

GSTより、いずれのノード、リンクとも高々 $\log n$ レベルにしか属さず、レベルの更新は単調増加なのでNCで使用される総メッセージ数は $\log n \times$

$O(n \log n) = O(n \log^2 n)$ となる。GSTで使われるメッセージ数は $O(m + n \log n)$ だから、アルゴリズムALEPのメッセージ複雑度は $O(m + n \log^2 n)$ 、メッセージサイズ $O(\log n)$ ビットとなる。

#### 4. 固有の識別子をもたず、故障を含むネットワークについて

これまでに固有の識別子をもたず、かつ故障を含むネットワークにおけるLEPについての結果は知られていない。本節では $G_n$ と $G_m$ について考察を行う。まず、 $G_n$ では、 $G_n$ 及び $G_r$ についての結果から解けないことは明らかである。次に $G_m$ においてもLEPが解けないことを示す。

まず、 $G_m$ においてLEPが条件(A)で解けるための条件を考える。あるプロセスが他のプロセスをリーダーとして承認を与えるということは、局所状態の差異による順序関係がつくことである。また、全体の過半数のプロセスからリーダーと承認を受けたプロセスが1つでも存在するとき、残りのプロセスは決して過半数から承認を得ることはない。各プロセスは $n$ を既知なので、結局 $G_m$ においてLEPが解けるための条件は全体の過半数のプロセス間で区別(存在の認識)がつくことである。

$G_m$ の非可解性を示す準備として $G_r$ においてLEPが解けないことを示す。図11のリング $G_2$  ( $\in G_r$ )を考える。リンクLは故障リンクとする。 $G_2$ 上でLEPを正しく解くアルゴリズムAが存在するとする。また、 $G_3$ は $G_2$ をLで切ったリング2つを対称的に配置し、一方は故障リンクで、もう一方は遅延の非常に大きなリンクで結合したものである。アルゴリズムAを $G_3$ に適用することを考える。Aが $G_2$ で条件(A)を満たすとすると、 $G_3$ の左右両部分でリーダーを選んでしまうので矛盾する。また、Aが $G_2$ において(A)を満たさず(B)を満たすとすると、Aを適用すると条件(A)で解けるような別のリングを切って $G_3$ の右半分に置き換えると、左半分は解けなかったと判断してしまい、矛盾する。よって $G_r$ においてLEPを解くことはできない。

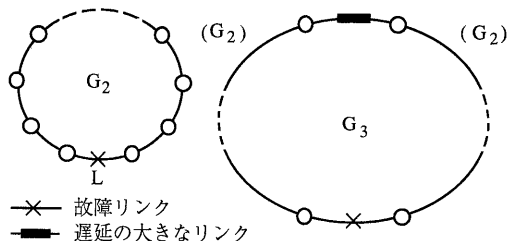


図11  $G_r$ におけるLEPの非可解性

あるネットワーク $G_4$  ( $\in G_{fin}$ ) を考える。各プロセスの識別子は固有ではないので、同じ識別子をもつ場合があり、そのときはポート数、ポートトラベリングなどの局所的な情報が全く同じプロセスを区別することは不可能である。LEPを解くためのあるアルゴリズム $A'$ が終了したとき、あるプロセス $P$ が識別できたプロセス集合からなる部分グラフを $G_d$ とする。 $G_d$ に含まれるプロセス数が全体の過半数を超えるような $G_d$ が存在すれば条件(A)で解ける。しかし、そうした $G_d$ が存在しない場合は各 $G_d$ で条件(B)を満たす必要がある。

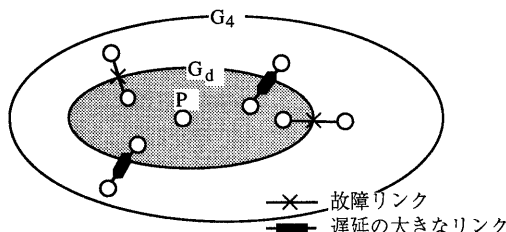


図12 識別可能な部分グラフ $G_d$

ここで図12のように考えると、 $G_d$ においてLEPを解くことは $G_4$ においてLEPを解くことと等価である。ところが $G_4$ においてはLEPを解くことは不可能であり、解けなかったことの検知はできない。結局 $G_d$  ( $\in G_{fin}$ ) においては条件(A)を満足しない場合が存在し、かつその場合条件(B)を満足できないので、 $G_{fin}$ においてLEPは解けない。

このことはランダム化アルゴリズムを用いても同じことが言える。ランダム化アルゴリズムとは各プロセスが独立に確率的選択を行えるアルゴリズムである。ここで独立な確率的選択とは識別子の変更のみに対して行われると仮定する。もし他の動作を確率的に行いたければ識別子に依存して決定的動作を行う手続きで代用できる。また識別子以外の値などを確率的に選択するとしても、識別子を元の識別子とその値の2項組にするなどすれば代用でき、この仮定は一般性を失わない。

ランダム化アルゴリズム $A''$ を用いて $G_{fin}$ においてLEPを解くことを考える。しかし、識別子変更を行っても新たに区別可能となるプロセスが増えるとは限らない。よって条件(A)を満たさない場合は存在する。また、 $G_4$ においてはランダム化アルゴリズムを用いてもLEPが解けないことは先の議論からも明らかである。従って $G_{fin}$ においてはランダム化アルゴリズムを用いてもLEPは解けない。

もちろんアルゴリズムをメッセージドリブンではなく、自発的に識別子を変更して送信するようなものに

すれば $G_{fin}$ においてもLEPは解ける。しかしメッセージ遅延は無限大にしない限りいくらでも大きく考えることができ、送信間隔は遅延に比べいくらでも小さくなる。したがってメッセージ数の上限は存在しないことになる。

## 5. むすび

本稿では故障を含むと仮定した任意形状の非同期ネットワークにおけるリーダー選挙問題について、各プロセスが固有の識別子をもつ場合とまたない場合について考察を行った。固有の識別子をもつ場合については故障リンク数の上限値が与えられたとした場合に対してメッセージ数 $O(m+n\log^2 n)$ のアルゴリズムの提案、またない場合については非可解性の証明を行った。

今後の課題としては提案アルゴリズムALEPのメッセージ複雑度の改善、また、今回はすべての故障がアルゴリズム開始時に起こっているとしたので形状が動的に変化するネットワークにおけるLEPの考察などが挙げられる。

最後に貴重な御意見を頂いた本学の山下雅史先生に深謝致します。本研究の成果の一部は文部省科学研究費補助金一般研究(B)(課題番号04452195)による。

## 文献

- [1] Afek, Y. and Saks, M.: "Detecting global termination conditions in the face of uncertainty," Proc. 6th Ann. ACM Symp. on Principles of Distributed Computing, pp.109-124 (1987).
- [2] Angluin, D.: "Local and global properties in networks of processors," Proc. 12th Ann. ACM Symp. on Theory of Computing, pp.82-93 (1980).
- [3] Bar-Yehuda, R., Kutten, S., Wolfstahl, Y. and Zaks, S.: "Making distributed spanning tree algorithms fault-resilient," Proc. 4th Symp. on Theoretical Aspects of Computer Science, pp.432-444, LNCS 247, Springer-Verlag (1987).
- [4] Goldreich, O. and Shrira, L.: "Electing a leader in a ring with link failures," Acta Informatica, Vol. 24, pp.79-91 (1987).
- [5] Itai, A. and Rodeh, M.: "Symmetry breaking in distributive networks," Proc. 22nd IEEE Symp. on Foundations of Computer Science, pp.150-158 (1981).
- [6] Lavalley, I. and Lavault, C.: "Spanning tree construction for nameless networks," 4th International Workshop on Distributed Algorithms, pp.41-56, LNCS 486, Springer-Verlag (1990).
- [7] Matias, Y. and Afek, Y.: "Simple and efficient election algorithms for anonymous networks," 3rd International Workshop on Distributed Algorithms, pp.183-194, LNCS 392, Springer-Verlag (1989).
- [8] Schieber, B. and Snir, M.: "Calling names on nameless networks," Proc. 8th Ann. ACM Symp. on Principles of Distributed Computing, pp.319-328 (1989).
- [9] Welch, J., Lamport, L. and Lynch, N.: "A lattice-structured proof technique applied to a minimum spanning tree algorithm," Proc. 7th Ann. ACM Symp. on Principles of Distributed Computing, pp.28-43 (1988).
- [10] 萩原, 増澤: "分散アルゴリズム," 情報処理学会誌, Vol.31, No.9, pp.1245-1256 (1990).
- [11] 亀田, 山下: "分散アルゴリズム," 出版予定.