

プロセッサとリンクの故障が混在する完全ネットワーク における分散アルゴリズム

西川直樹†

増澤利光††

都倉信樹†

†大阪大学基礎工学部

††大阪大学情報処理教育センター

プロセッサとリンクの故障が混在する完全ネットワークにおける放送問題とリーダー選択問題の通信計算量について考察し、リンクとプロセッサがそれぞれ高々 f_L, f_P 個故障する完全ネットワーク (n はプロセッサ数) で以下の結果を得た。

(1) 方向感覚付き完全ネットワーク上の放送問題は $\Theta(n + f_L \cdot f_P + f_L \cdot \log f_L)$

(2) 方向感覚のない完全ネットワーク上の放送問題は $\Theta(n + n \cdot f_L)$

(3) 方向感覚付き完全ネットワーク上のリーダー選択問題は $\Theta(n + k(f_P + f_L) + f_L \cdot f_P + f_L \cdot \log f_L)$ 。ただし、 k は起動プロセッサ数を示す。

Distributed Algorithms in Complete Networks with Link and Processor Failures

Naoki Nishikawa†

Toshimitsu Masuzawa††

Nobuki Tokura†

† Faculty of Engineering Science
Osaka University

†† Education Center for Information Processing
Osaka University

This paper considers fault-tolerant distributed algorithms in asynchronous complete networks with undetectable fail-stop link and processor failures. The problems studied in this paper are the broadcast problem and the leader election problem, and the followings are shown for complete networks with n processors and at most f_L faulty links and at most f_P faulty processors.

(1) The message complexity of broadcast is $\Theta(n + f_L \cdot f_P + f_L \cdot \log f_L)$, if global sense of direction is available, while

(2) the message complexity is $\Theta(n + n \cdot f_L)$, if no global sense of direction is available.

(3) The message complexity of leader election is $\Theta(n + k(f_P + f_L) + f_L \cdot f_P + f_L \cdot \log f_L)$ where k represents the number of initiator processors start the algorithm, if global sense of direction is available.

1. まえがき

通信用リンクで結合されている複数台のプロセッサ(すべてのプロセッサは対等であり、ホストプロセッサなるものはない。以下プロセッサをPEと略記)に、ある問題を解くために必要な情報が分散している状況で、それらの情報を交換しながらその問題を解くアルゴリズムは分散アルゴリズム(distributed algorithm)と呼ばれる。(但し、PE間には共有メモリはなく、PE間通信はリンクを通じたメッセージ交換によってのみ行なわれる。)分散アルゴリズムは非同期式ネットワーク(リンクを伝えるメッセージの伝送遅延、PE間の動作速度の比に関して、有限であること以外には何も仮定しないネットワーク)で考えられることが多い。これまでにリーダー選択問題や生成木構成問題などを解く多くの分散アルゴリズムが提案されている(5), (7), (9), (11)-(13)。

リンクやPEの故障が予想されるときに、それらの故障にもかかわらず種々の問題を解く分散アルゴリズムについても研究されている。特に非同期ネットワークではリンクやPEの故障を検知できないので、故障耐性のあるアルゴリズムを開発することは容易ではない。これまでに提案された故障耐性のあるアルゴリズムのほとんどは、リンクまたはPEどちらか一方のみの故障に限定して考えられている(1)-(4), (6), (8), (10)。それらのアルゴリズムは、その解法をPEまたはリンクには故障がないという仮定に強く依存していることが多く、PEとリンクの故障が混在しているような状況では問題を解くことができない。そこで、本稿ではリンクとPEの故障が混在している場合の分散アルゴリズムについて考察する。

分散アルゴリズムに於いては通信計算量(最悪時評価)やその他の評価基準をより改善するために対象となるネットワークの形状にトポロジ制限を与えることがある。このような目的でリングや完全ネットワークが用いられるがここでは方向感覚付き完全ネットワーク(13)を用いる。

本稿では、放送問題(ある一台のPEのもつ情報を故障していない全PEに知らせる問題)について考察する。分散アルゴリズムにおいて、多くの場合ネットワーク中の全PEが結果あるいは終了の事実を認識してアルゴリズムを終了することから、他の問題の一部として用いることができる基本的な問題である。放送問題を解く通信計算量については、方向感覚付き完全ネットワークでPE故障のみの場合は $\Theta(n)^\dagger$ 、リンク故障(高々 f_L 個)のみの場合は $\Theta(n+f_L \log f_L)$ であり、方向感覚のない完全ネットワークでは、PE故障のみの場合は $\Theta(n)$ 、リンク故障(高々 f_L 個)のみの場合は $\Theta(n+n \cdot f_L)$ であることが知られている(10)。

本稿では、高々 f_P 個のPEと f_L 個のリンクに故障がある(ただし $f_P+f_L \leq n-2$)場合に、方向感覚付き完全ネットワークにおける放送問題を解く通信計算量 $O(n+f_P \cdot f_L + f_L \cdot \log f_L)$ の分散アルゴリズムと、方向感覚のない完全ネットワーク上で、放送問題を解く通信計算量 $O(n+n \cdot f_L)$ の分散アルゴリズムを示す。そして、これ

\dagger 実数から実数への関数 $f(n), g(n)$ に対し、正定数 c, n_0 が存在し、すべての $n \geq n_0$ に対し、 $f(n) \leq cg(n)$ なるとき、 $f(n) = O(g(n))$ と書く。

同様に、 $f(n) \geq c'g(n)$ なるとき $f(n) = \Omega(g(n))$ と書く。

$f(n) = O(g(n))$ かつ、 $f(n) = \Omega(g(n))$ のとき、 $f(n) = \Theta(g(n))$ と書く。

らのアルゴリズムの通信計算量がどちらもオーダー的に最適であることを示す。

さらにこのアルゴリズムを用いて、方向感覚付き完全ネットワークでリーダー選択問題を解く通信計算量 $O(n+k(f_P+f_L)+f_P \cdot f_L + f_L \cdot \log f_L)$ (k は起動プロセッサ数)のアルゴリズムを示す。また、このアルゴリズムの通信計算量がオーダー的に最適であることを示す。

(なお、文献(11)は方向感覚付き完全ネットワークにおけるリーダー選択問題の通信計算量についてPE故障(高々 f_P 個)のみの場合は $\Theta(n+k \cdot f_P)$ 、リンク故障(高々 f_L 個)のみの場合は $\Theta(n+k \cdot f_L + f_L \cdot \log f_L)$ であることが示されているので本稿の結果はその拡張とみることができる。)

2. 諸定義

2.1 ネットワーク

[定義1] ネットワーク N は二項組 $N=(P, L)$ で定義され、以下の性質を持つ。

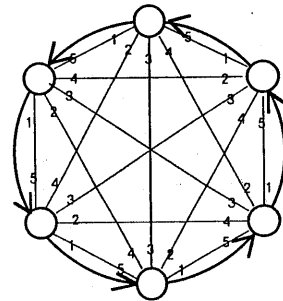
- (1) プロセッサ(以下では、PEと略記する)間の通信はリンクを介したメッセージ交換によってのみ行える。
- (2) リンクは両方向の通信が独立に(反対方向の通信の影響を受けることなく)行える双方向リンクである。
- (3) PEは互いに異なる識別子を持つ。PE u の識別子を $id(u)$ と表す。
- (4) ネットワークは非同期式である。つまり、メッセージがリンクを伝わる伝送遅延は有限だが可変であり前もって分からない。
- (5) 各PE (u とする)は、 $deg(u)$ 個($deg(u)$ は u に接続するリンクの数を表す)のポート $\{a \mid 1 \leq a \leq deg(u)\}$ を持つ。各ポートは u に接続する各リンクと1対1に対応する。 u のポート a がリンク (u, v) に対応するとき、 $v = u[a]$ と表す。ただし、 $u[0]$ は u 自身を表す。
- (6) リンクはFIFO (First-In First-Out)とする。つまりリンク (u, v) を用いてPE u が v に送信したメッセージは u が送信した順に v が受信する。□

以下、本稿ではネットワーク全体のPE数を n で表す。

2.2 完全ネットワークと方向感覚

[定義2] 任意の相異なるPE間にリンクが存在するネットワークを完全ネットワーク(以後、CNと略記)という。CN N が次の(条件SD)を満足するとき、 N は方向感覚付き完全ネットワーク(Complete Network with a Sense of Direction, 以後、CN-SDと略記)である(図1)。

なお、任意の整数値 x, y に対して式 $(\overline{x-y})$ は $x-y$ の法 n の非負剰余を示す。



→:ハミルトン有向閉路

図1:方向感覚付き完全ネットワーク

(条件SD) u を任意のPEとし, a, b を $(1 \leq a, b \leq n-1)$ を満たす任意の正整数とする. $v = u[a], w = u[b]$ とすると, $w = v[b-a]$ が成立する. (図2) □

つまり, CNの1つの有向ハミルトン閉路 H を考え, H 上で距離が a のPEとはポート a を介して隣接しているCNがCN-SDである.

2. 3 分散アルゴリズム

[定義3] 分散アルゴリズム (以下では単にアルゴリズムという) は以下の性質を持つ.

(1) 各PEが実行するプログラムは, 受信したメッセージに対する処理 (どのような内部計算をし, どのようなメッセージをどのポートから送信するか) を表す.

(2) すべてのPEのプログラムは同一である. 但し, 各PE u のプログラムでは, $id(u)$, ポート数 $deg(u)$, f_p と f_L (ネットワーク全体での故障PEと故障リンクそれぞれの数の上限) を利用できる.

(3) いくつかのPEは, 最初に他のPEからメッセージを受信しなくとも自主的にメッセージを送信する. このようなPEを起動PEと呼ぶ. 起動PE以外のPEは, 他のPEからメッセージを受信してからプログラムの実行を開始する. □

2. 4 放送問題とリーダー選択問題

[定義4] $N=(P, L)$ を任意のネットワークとする. N における放送問題とは, あるPE $u(\in P)$, 放送PEと呼ぶ)の持つあるメッセージ M を N 中の全非故障PEに伝達する問題である. ただし, 放送PEのみが起動PEになるとする. □

[定義5] $N=(P, L)$ を任意のネットワークとする. N におけるリーダー選択問題とは, あるPE $v(\in P)$ をリーダーとして選択する問題である. 具体的には v が「自分がリーダーである」ことを知り, v 以外のPEが「リーダーの識別子は $id(v)$ である」ことを知る問題である. ただし, 1個以上の任意のPEが起動プロセッサになるとする. □

[定義6] A を任意のアルゴリズム, $N=(P, L)$ を任意のネットワーク, Q をある分散問題としたときに, リンクを伝搬中のメッセージが存在しない状況で, 各起動PEが A の実行を開始したとする. このとき起動PEの起動のタイミングや各PEの動作速度, 伝送遅延などの違いにより, 様々な実行経過が考えられる. これらのどの実行経過においても, すべての非故障PEが有限時間内に停止し, すべての非故障PEが停止したとき問題 Q が解かれているとき, 「 A は N における問題 Q を解く」という. □

2. 5 通信計算量

ここでは分散問題を解くために要するメッセージの個数に関する概念を定義する.

[定義7] アルゴリズムを実行したときにネットワーク全体で送信されるメッセージの総数 (故障PEや故障リンクに対して送信されたメッセージを含む) を, そのアルゴリズムの「通信計算量」とよぶ. 但し, 起動PE,

故障PEやリンクの遅延などの違いによりさまざまな実行経過があるが, その中で最も多くメッセージを用いる場合のメッセージ数で評価する (最悪時評価). □

1つのメッセージで任意に多くのデータを通信できると考えるのは現実的ではないので, 次の仮定をおく.

[仮定1] 1つのメッセージ長は $O(\log n)$ ビットとする. また, 各PEの識別子は $O(\log n)$ ビットで表せるとする. □

[仮定2] 放送問題において M (全PEに伝達しなければならないメッセージ) は $O(\log n)$ ビットで表せるとする. □

2. 6 故障

[定義8] PEに故障があると, そのPEは送受信いずれも行えなくなる. また, そのPEが起動することはない. リンクに故障があると, そのリンクを用いた両方向の通信ができなくなる. □

以後, 故障PEのことをFPE, 故障していないPEをNPEと呼ぶ.

[仮定3] FPEおよび故障リンクを検知するような能力はネットワークにはない. □

[仮定4] アルゴリズム実行中はPEやリンクの故障が新たに生じたり, FPEや故障リンクが復旧したりしない. □

また, ネットワークにおいて, FPEと故障リンクの数には上限があり, それをそれぞれ f_p と f_L とする. 各PEはこの値を利用してよいものとする.

[定義9] $N=(P, L)$ に対して $Live(N)=(P', L')$ は, 全てのNPEとリンクよりなる部分グラフである. つまり $P' = P - F_p$, $L' = \{e \in P' \times P' | e \in L - L_f\}$ とする. □

2. 7 グラフ

[定義10] グラフ G は, 頂点の集合 V と V の相異なる二要素の二項組の集合 E の組 (V, E) で定義される.

グラフ $G=(V, E)$ において, 各 $i(1 \leq i < m)$ に対し $(u_i, u_{i+1}) \in E$ を満たす相異なる頂点の系列 $\langle u_1, u_2, \dots, u_m \rangle$ を $u_1 - u_m$ 道と呼び, 道に含まれる辺の数 $m-1$ をその道の長さと呼ぶ. 2頂点 u, v の間に長さ2以上の $u-v$ 道が存在し, $(u, v) \in E$ であるとき, これを閉路という. 任意の2頂点に対し, その長さが最小である $u-v$ 道を $u-v$ 最短道と呼ぶ. この $u-v$ 最短道の長さを $u-v$ 距離と呼び, $d_G(u, v)$ で表す.

グラフ $G=(V, E)$ が閉路を含まなければ G を木と呼ぶ. 木のある頂点 u と任意の頂点 v の距離が d 以下であるとき, 木 G は u を根とする高さ d の木であるという.

グラフ $G=(V, E)$, $G'=(V', E')$ に対して $V' \subseteq V, E' \subseteq E$ がなりたつとき, G' を G の部分グラフという. 特に, $V=V'$ かつ G' が木であるとき G' を G の生成木と呼ぶ. □

3. CN-SDでのプロセッサとリンクの故障が混在する場合の放送問題

FPEが高々 f_p 個, 故障リンクが高々 f_L 個存在するCN-SDで, 放送問題を解くプロトコル *deforest* ("森林開拓"の意) を示す. ただし, $f_p + f_L \leq n-2$ とする. この条件は $Live(N)$ が連結であることを保証する.

3. 1 プロトコル *deforest*

放送プロトコル *deforest* は3つのステージからなる. (ステージ1) 最初に放送PE u から全てのPEに対して直接メッセージ M を送信する. M を受信したPEは M を受け取ったことを u に返信する. f_p 個のFPEが存在し f_L 個の故障リンク全てが u に接続している可能性があるので, u は $f_L + f_p$ 個のPEからの返信は期待できない. $n - f_L - f_p - 1$ 個の返信を受信すれば u はステージ2に移行する.

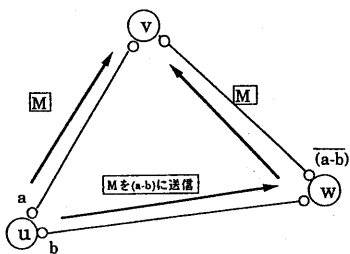


図2: 方向感覚を利用した位置情報の伝達例

(ステージ2) 放送PE u が M を (ステージ1で M を受信しなかった) PE v に伝達するには、リンク (u, v) が故障している場合、他のPEを介して中継する必要がある。もちろんこの中継PEと v との間のリンクも故障している可能性があるため、複数の中継PEを作る必要がある(図3)。あるPE w を中継PEとして指定したとする、 u はまだ返信を受け取っていないPEの位置を方向感覚を用いて w に連絡し、 w はそれらのPEに M を送信する(図2)。 w に対して M を受信したことに対する返信がきたならば、 w はそのことを u に連絡する。各PEに対して u と w の両方から M を送信するので、故障リンクによって M を受信できないNPEの数は半分になる。従って、返信が期待できないPEの数は、

$\lfloor f_L/2 \rfloor + f_p$ 個になる。 M は (ステージ1に受信したのもも含めて) $n - \lfloor f_L/2 \rfloor - f_p - 1$ 個の返信を受信したならば、中継PEの数を1つ増やす。このようにして次々に中継PEを増やしていく。 m 個の中継PEを利用すれば、 $n - \lfloor f_L/(m+1) \rfloor - f_p - 1$ 個のPEが M を受信することが保証される。従って f_L 個の中継PEを用いれば f_p 個以外のプロセッサに M を伝えることができる。

(ステージ3) 高々 f_p 個のFPEが存在しうるので、 f_p 個のPEに対しては (故障していて永久に返信がこないかも知れないので、) M の伝達の確認を待つことができない。そこでこれら f_p 個のPEに対しては $f_L + 1$ 個の隣接PEから M (とアルゴリズムの停止命令) を送信する。これにより、それらのPEがNPEならば故障リンクに妨げられることなく、有限時間以内に M を受信 (してその後停止) する。

なお、各ステージにおける動作を森林開拓と見立てて (最初「woods(森林: M を受信していないPE)」が生い茂っていると考え、「flat(平地: M を受信したPE)」が広がっていく) アルゴリズムの各名称をつけた。

3. 1. 1 PEの状態

各NPE v のとりうる状態は以下の6つである。

woods: M を受信していないPEの状態、放送PE u 以外のPEの初期状態でもある。

c_flat: u から直接 M を受信したPE (リンク (u, v) が故障していないこと判っている) の状態 (close-flat: “近くの平地”の略)。

f_flat: 中継PEからのみ M を受信したPE (リンク (u, v) が故障しているかどうか不明) のPEの状態 (far-flat: “遠くの平地”の略)。

colony: 中継PEの状態 (“開拓村”の意味)。

col_bs: 放送PEの状態 (colony-base: “開拓起源地”の略)。

halt: 停止状態。

なお、以後 “状態XXXのPE” を単に “XXX” と記述する。

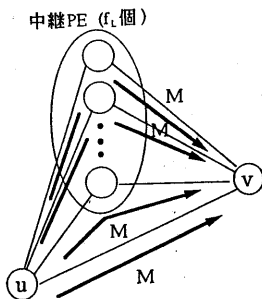


図3: リンク故障を回避するために中継PEを用いる

3. 1. 2 PEの局所変数

各プロセッサが用いる局所変数は以下の通り。

・全NPEが用いる局所変数

status: PEの状態を表す変数。

dir_cb (direction of colony-base): 放送PEに通じる (非故障リンクよりなる) 道へのポート番号を記録する。

・中継PE (colony) で用いる局所変数

f-table [1..n-1]: 自分 (v とする) と非故障リンクでつながっていることの保証された、状態が f_flat のPEを記録する。 $f_table[t] = true$ ならば、 $v[t]$ の状態が f_flat であることを示す。初期値はすべて $false$ 。

・放送PE u (col_bs) で有効な局所変数

netdata [1..n-1] (network data): 他の全PEの状態を記録する。 $netdata[a]$ は $u[a]$ の状態を記録する。初期値はすべて $woods$ 。

#woods (number of woods): $netdata$ 中で状態が $woods$ のPEの総数。

#colony (number of colony): $netdata$ 中で状態が $colony$ のPEの数を示す。

#c_flat (number of c_flat): $netdata$ 中で状態が c_flat のPEの数を示す。

netflat [1..n-1] (network f_flat): 状態が f_flat のPE $u[a]$ が、 $colony$ $u[b]$ を介して放送PE u に連結しているとき、 $netflat[a] = b$ が成り立つ。初期値はすべて0。

3. 1. 3 メッセージ

以下の9種類のメッセージを用いる。

《**fell1**(M)»: 最初に放送PE u から送信される (fellは “伐採” の意味)。

《**clear1**》»: 《**fell1**》に対する返信 (clearは “伐採完了” の意味)。

《**colonize**(d)»: u から $colony$ 、または c_flat (v とする) に送信される。 v はこのメッセージを受信すると、

$v[d]$ に 《**fell2**》を送り、さらに v が c_flat ならば、 v は $colony$ になる。 (colonizeは “開拓村設立” の意味)

《**fell2**(M)»: $colony$ から送信される。

《**clear2**》»: 《**fell2**》に対する返信。

《**clear3**(d)»: $colony$ が、 $u[d]$ から 《**clear2**》を受信したことを u に知らせるメッセージ。

《**connect**》»: f_flat v が、リンク (u, v) が非故障であることを u に連絡するメッセージ (connectは “開拓起源地との連絡道完成” の意味)。

《**order1**(d)»: ステージ3に放送PE、 $colony$ より送信する。受信したPE v は $v[d]$ に 《**burn**》を送信する (orderは “命令” の意味)。

《**order2**(d, d')»: ステージ3に u が $colony$ v に送信し、このメッセージを受信した v は $v[d]$ に 《**order1**(d')》を送信する。

《**burn**(M)»: (故障している可能性のある) まだ u がメッセージ M を受信したという返信を受信していないPEに対して送信され、かつそのPEを受信後に停止させる (burnは “切り開いた土地を農地にするために) 焼き払う行為” の意味)。

《**halt**》»: 停止命令。

3. 1. 4 PEの動作

簡単のため、アルゴリズムの記述に以下の送信命令を用いる。なお、これらの命令を実行するPEを v とする。

・ **send** ($i, message$): $v[i]$ に ($message$) を送信する。

・ **r-rec** ($t, message$): v はメッセージを受信するまで受信命令を繰り返す。最初に受信したメッセージは $message$ に代入され、受信したポートは t に代入される。

col_bs uは#woodsの値によって、次のように3つのステージを実行する。(u以外のPEにはステージが関係なく、statusと受信メッセージによって動作する。)

ステージ1 [deforestation (荒い伐採)] :
 #woods>fp+f_Lの間。
 ステージ2 [uproot onebyone (一本ずつ引き抜く)]
 fp+f_L≥#woods>fpの間。
 ステージ3 [burn up (残りを燃やしてしまう)]
 fp≥#woodsの間。

放送PEの動作

```
(ステージ1 : [deforestation])
status:=col_bs;
#woods:=n-1; #colony:=0; #c_flat:=0;
for i=1 to n-1 do
  netdata[i]:=woods;
  netflat[i]:=0;
end for
for i=1 to n-1 do
  send(i, «fell1(M)»);
end for
while #woods>fp+fL do
  r-rec(t,message);
  if message= «clear1» then
    #woods:=#woods-1;
    netdata[t]:=c_flat;
    #c_flat:=#c_flat+1;
  end if
end while
(ステージ2 : [uproot onebyone])
while fp<#woods do
  if #woods≤fp+⌊fL/(#colony+1)⌋
  and #c_flat≥1 then
    netdata[t]=c_flatなるtを1つ選択;
    netdata[t]:=colony;
    #colony:=#colony+1;
    #c_flat:=#c_flat-1;
    for i=1 to n-1 do
      if netdata[i]=woods then
        send(t, «colonize(i-t)»);
      end if
    end for
  end if
  r-rec(t,message);
  case message of
    «clear1»:
      #woods:=#woods-1;
      netdata[t]:=c_flat;
      #c_flat:=#c_flat+1;
    «clear3(d)»:
      if netdata[d]=woods then
        #woods:=#woods-1;
        netdata[d]:=f_flat;
        netflat[d]:=t;
      end if
    «connect»:
      if netdata[t]=woods then
        #woods:=#woods-1;
      end if
      netdata[t]:=c_flat;
      #c_flat:=#c_flat+1;
```

```
end case
end while
(ステージ3 : [burn up])
netdataを参照して、colony,c_flat,f_flatの中から順にfL個のPEを選ぶ。つまり、(1)#colony≥fLならばcolonyだけから選び、(2)#colony<fLかつ#c_flat+#colony≥fLならば、fL-#colony個のc_flatとcolony全部を選び、(3)それ以外の場合には、fL-#colony-#c_flat個のf_flatとcolony,c_flatのそれぞれ全部を選ぶ。
選んだfL個の各PE (u[d]とする) に対して次の動作を実行する。
```

```
for i=1 to n-1 do
  if netdata[i]=woods then
    if netdata[d]=c_flat
    or netdata[d]=colony then
      send(d, «order1(i-d)»);
    else
      d':=netflat[d];
      send(d', «order2(d-d',i-d)»);
    end if
  end if
end for
さらに以下の動作を行なう。
for i=1 to n-1 do
  if netdata[i]=c_flat
  or netdata[i]=colony then
    send(i, «halt»);
  end if
end for
```

放送PE以外のPEの動作

```
イ) status=woodsのPEの動作
while status=woods do
  r-rec(t,message);
  case message of
    «fell1(M)»:
      status:=c_flat;
      dir_cb:=t;
      send(t, «clear1»);
    «fell2(M)»:
      status:=f_flat;
      dir_cb:=t;
      send(t, «clear2»);
    «burn(M)»:
      dir_cb:=t;
      status:=halt;
  end case
end while
ロ) status=c_flatのPEの動作
while status=c_flat do
  r-rec(t,message);
  case message of
    «colonize(d)»:
      status:=colony
      send(d, «fell2(M)»);
    «order1(d)»:
      send(d, «burn(M)»);
    «halt»:
      status:=halt;
```

```

end case
end while
ハ) status=f_flatのPEの動作
while status=f_flat do
  r-rec(t,message);
  case message of
    《fell1(M)》:
      status:=c_flat;
      dir_cb:=t;
      send(t,《connect》);
    《order1(d)》:
      send(d,《burn(M)》);
    《halt》:
      status:=halt;
  end case
end while
ニ) status=colonyのPEの動作.
while status=colony do
  r-rec(t,message)
  case message of
    《colonize(d)》:
      send(d,《fell2(M)》);
    《clear2》:
      f-table[t]:=true;
      send(dir_cb,《clear3(t-dir_cb)》);
    《order1(d)》:
      send(d,《burn(M)》);
    《order2(d,d')》:
      send(d,《order1(d')》);
    《halt》:
      for i=1 to n-1 do
        if f-table[i]=true then
          send(i,《halt》);
        end if
      end for
      status:=halt;
    end case
end while

```

3. 2 アルゴリズムの正当性

[定理1] アルゴリズムdeforestは、有限時間内に放送問題を解く。

(証明) 最初に放送PE u が有限時間内に停止することを示す。

ステージ1は u が《fell1》を他の全てのPEに送信し、 $n-f_p-f_L-1$ 個の《clear1》が返信されるのを待つだけなので有限時間内に終了する。ステージ3も u が有限個数のメッセージを一方向的に送信するだけなので、有限時間内に終了する。

ステージ2が有限時間内に停止しないと仮定する。 u が送信する有限個の《fell1》と《colonize》に対する有限個のメッセージの送返信によってステージ2は動作しているので、いずれ以下のような状態になる。

「(放送PE u とcolonyからwoodsへのリンクがすべて故障して) woodsから c_flat や f_flat になるものがなく、 c_flat がいずれもcolonyにならず、どの f_flat も c_flat とならない」

上記の状態が発生したと仮定する。もし c_flat が一つでも存在すると、 u とcolonyから残りのwoodsへのリンクが全て故障しているので、これらのリンク数は f_L を越え

ず、

$$\#woods \leq f_p + \lfloor f_L / (\#colony + 1) \rfloor$$

が成り立っている。すると、 u はその c_flat へ

《colonize》を送信して有限時間内に c_flat はcolonyになる。よって、 c_flat は一つもないことになる。

すると $f_p + f_L \leq n-2$ の条件からネットワークはNPEが非故障リンクで連結であるので、 u から非故障リンクを少なくとも三つ以上介さないと残りのwoodsに対してつながらないことになる。しかし、これが成り立たないことを以下の補題1と補題2で示す。

[補題1] 頂点数 n の完全グラフ $G=(V,E)$ の任意の部分グラフ $G'=(V,E')$ について、 $|E|-|E'| \leq n-2$ ならば、任意の頂点を根とする高さ2以下の生成木が G' に存在する。

(証明) $|E|=|E'|$ のときはあきらか。 $|E| > |E'|$ のとき、 $E-E'=F$ としたときに $(u,v) \in F$ なる $u,v \in V$ に対し、 $dg_G(u,v) > 2$ とする。このとき u,v 以外の任意の頂点 $w \in V$ について $(u,w) \in F$ または $(v,w) \in F$ となる。このような w のとりかたが $n-2$ 通りあり、 $(u,v) \in F$ なので、 $|E|-|E'| \geq n-1$ となって矛盾する。よって、任意の u,v について、 $dg_G(u,w) \leq 2$ が成り立つ。よって、任意の頂点を根とする高さ2以下の G' の生成木が存在する。□

[補題2] PE数 n 、故障リンク数が高々 f_L 、FPE数が高々 f_p のCN(ただし $f_p + f_L \leq n-2$)において、任意のNPE u に対して「 u を根とする高さ2以下の生成木」が少なくとも一つは存在する。

(証明) FPEとそれに連結するリンクをすべて除外してできるネットワークは頂点数 $n-f_p$ のCNである。また $f_L \leq n-f_p-2$ である。よって、補題1より高さ2以下の生成木は存在する。□

よって先の仮定は否定されステージ2も有限時間内に終了する。

故に各ステージは有限時間内に終了し、 u は有限時間内に停止する。

次に u 以外のNPEが M を受信し、有限時間内に停止することを示す。

ステージ3までに放送PEがメッセージ M の受信を確認しているNPE(u のnetdataでの状態がcolony, c_flat , f_flat のいずれかのもの)は、途中で状態が f_flat から c_flat に変化したPEを除いて、 M を受信したのと同じ道筋で《halt》を受信する。途中で状態が f_flat から c_flat に変化したPEは、 u がステージ3を開始する前に《connect》を受信されていれば u から、そうでなければcolonyから《halt》を受信する。

ステージ3までに放送PEがメッセージ M の受信を確認していなかったNPEに対して、ステージ3で f_L+1 個の異なるプロセッサから《burn》を送信するので、少なくともそのうちの1個は受信される。

これらのことから全てのNPEは有限時間内に M を受け取って、停止する。

故にアルゴリズムdeforestは有限時間内に放送問題を解く。□

3. 3 アルゴリズムの評価

[補題3] アルゴリズムdeforestの実行中にcolonyとなるPEは高々 f_L 個である。

(証明) col_bs はステージ2においてのみ新しいcolonyを作ろうと《colonize》を送信する。このとき、

$$\#woods \leq f_p + \lfloor f_L / (\#colony + 1) \rfloor$$

が成り立つ。また、まだステージ3に移行していないので $f_p < \#woods$ が成立つ。従って、

$$\#colony < f_L$$

を導くことができる。□

[補題4] アルゴリズムdeforestの実行中に送信される《colonize》は高々 $fp \cdot f_L + O(f_L \log f_L)$ 個である。

(証明)

PEがcolonyとなったときに《colonize》はそのときの#woodsの値だけ送信される。PEがcolonyとなるとき

$$\#woods \leq fp + \lfloor f_L / (\#colony + 1) \rfloor$$

が成り立つので、k番目にcolonyとなるPEには高々 $(fp + (f_L/k))$ 個の《colonize》が送信される。また、補題4よりcolonyとなるPEの数は高々 f_L 個であるので、アルゴリズム全体の《colonize》の送信回数は高々

$$f_L \sum_{k=1} (fp + (f_L/k)) = fp \cdot f_L + f_L \log f_L + O(f_L)$$

となる。□

[定理2] アルゴリズムdeforestの最悪時通信計算量は $O(n + fp \cdot f_L + f_L \log f_L)$ である。

(証明) 《fell1》はステージ1に放送PE uが他の全てのPEに送信するだけなのでその総数は $n-1$ 個である。

《clear1》はその返信なので高々 $n-1$ 個である。

ステージ2において、《colonize》が1つ送信されると、《fell2》が1つ送信され、その返信として

《clear2》、さらにそれを受信したcolonyによって《clear3》が送信される。よって《fell2》、《clear2》、《clear3》の送信総数は《colonize》の送信個数を越えない。補題4よりこれらのメッセージの送信個数はいずれも高々 $fp \cdot f_L + O(f_L \log f_L)$ 個である。

ステージ3の動作より、《order1》は選ばれた f_L 個の中継PEへそれぞれ fp 個ずつ送信される。従って送信個数は高々 $fp \cdot f_L$ 個である。《order2》は

《order1》の送信個数を越えることはなく、《burn》は《order1》の送信個数と同じで、どちらの送信個数も高々 $fp \cdot f_L$ 個である。また、《halt》の送信個数は高々 $2n$ である。

よってアルゴリズム全体での最悪時通信計算量は $O(n + fp \cdot f_L + f_L \log f_L)$ である。□

3. 3 CN-SDでの放送問題の通信計算量の下界

FPEと故障リンクの混在するCN-SDでの放送問題を解くアルゴリズムdeforestの通信計算量はオーダー的に最適であることを示す。

[補題5] ⁽¹⁰⁾ PE数 n 個のCN-SDに f_L 個($\leq n-2$)のリンク故障のある場合の放送問題の通信計算量の下界は $\Omega(n + f_L \log f_L)$ である。

[定理3] PE数 n 個のCN-SDに f_L 個のリンク故障と fp 個のPE故障がある場合の放送問題の通信計算量の下界は $\Omega(n + fp \cdot f_L + f_L \log f_L)$ である。ただし $fp + f_L \leq n-2$ とする。

(証明) $fp = 0$ の場合を考えれば、補題5より、下界 $\Omega(n + f_L \log f_L)$ が証明できる。

放送問題を解く任意のアルゴリズムをAとする。

ネットワーク $N=(P, L)$ に対して、 P' を $|P'| = fp+1$ かつ $u \notin P'$ を満たすPの任意の部分集合とし、

$L' = \{(v, w) \mid v \in P', w \in P - P'\}$ とする。Aは、 L' のうち少なくとも $(fp+1) \cdot (f_L+1)$ 個のリンクにメッセージを送信するような実行経過があることを示す。

Aが L' に $(fp+1) \cdot (f_L+1)$ 個より少ないリンクにメッセージを送信して終了すると仮定すると、 $|L'| \leq f_L$ (ただし、 $L' = \{(x, w) \in L' \mid (x, w) \text{はAがメッセージを送信したリンク}\}$ が成り立つ) ようなあるNPE $x(x \in P')$ が存在する。すると、 P' の x 以外のPEが全てFPEで、 L' が全て故

障リンクである場合、 x はMを受信することができない。よって先の仮定は誤っている。

故にAは少なくとも $(fp+1) \cdot (f_L+1)$ 個のメッセージを送信し、放送問題の通信計算量の下界は $\Omega(n + fp \cdot f_L + f_L \log f_L)$ である。□

4. 方向感覚のないCNにおける放送問題

[補題6] ⁽¹⁰⁾ PE数 n 個のCNに f_L 個($f_L \leq n-2$)のリンク故障のある場合の放送問題の通信計算量の下界は $\Omega(n + n \cdot f_L)$ である。

[定理4] PE数 n 個のCNに f_L 個のリンク故障と fp 個のPE故障がある場合(ただし $fp + f_L \leq n-2$)の放送問題の通信計算量は $\Theta(n + n \cdot f_L)$ である。

(証明) 3. で述べたdeforestアルゴリズムを次のように変更する。(変数などの細かい変更は省略)

(1)新たにcolonyとなったPEは放送PE uに通じる以外全てのポートへ《fell2》を送信する。

(2)《order1》を受信したPEは《order1》を受信したポート以外の全てのポートへ《burn》を送信する。

(3)《order2》はidによって中継PEを指定する。

このような変更を加えてできるプロトコルは、《fell2》と《burn》の送信個数のみがdeforestよりも多くなる可能性がある。

《fell2》がcolonyごとに $(n-2)$ 個送信され、(依然としてcolonyの総数が高々 f_L 個なので)《fell2》はアルゴリズム中に高々 $n \cdot f_L$ 個送信される。《burn》はステージ3で選ばれた f_L 個のPEから高々 $(n-2)$ 個ずつ送信されるので、全部で高々 $n \cdot f_L$ 個送信される。よってこのアルゴリズム全体での通信計算量は $O(n + n \cdot f_L)$ となる。

補題6より $\Omega(n + n \cdot f_L)$ が成り立つので、CNにおける放送問題の通信計算量は $\Theta(n + n \cdot f_L)$ である。□

定理2、4から方向感覚が放送問題の通信計算量を減らす上で有用であることが判る。

5. 故障の混在するモデルにおけるリーダー選択問題

ここでは $fp + f_L \leq n-2$ かつ $fp < n/2$ ($fp \geq n/2$ では、故障リンクがなくてもリーダー選択問題を解くアルゴリズムは存在しない⁽⁸⁾)のときのリーダー選択問題を解くアルゴリズムdfLEAD ($fp + f_L < n/4$ のときに用いるdfLEAD1と、 $n-2 \geq fp + f_L \geq n/4$ のときに用いるdfLEAD2よりなる)を示す。

5. 1 アルゴリズムdfLEAD1

文献(10)のアルゴリズムLEAD(f) ($fp=0$, $f_L < n/4$ のときのリーダー選択アルゴリズム)は、 $fp + f_L < n/4$ の場合でも、1つのPEをリーダーとして選択する部分に関しては正しく動作する。(この部分の通信計算量は $O(n + k \cdot (fp + f_L))$ である)リーダー選択後、リーダーの識別子を放送する通知フェーズをdeforestプロトコルで行なえば故障の混在するモデルのリーダー選択問題を解くことができる。これがdfLEAD1であり、全体の通信計算量は $O(n + k \cdot (fp + f_L) + fp \cdot f_L + f_L \cdot \log f_L)$ になる。

5. 2 アルゴリズムdfLEAD2

プロトコルdeforestでは、中継PEを最大 f_L 個まで作ることによって全NPEにメッセージを送信した。deforest実行時において、中継PEを m 個用いるとMを受けとれないNPEの個数は高々 $fp + \lfloor f_L / (m+1) \rfloor$ 個である。従って中継PEを2つ作れば、 $(n - fp - \lfloor f_L / 3 \rfloor - 1)$ 個のNPEには確実にメッセージを伝達できる。一方、 $fp + f_L \leq n-2$ と $fp < n/2$ より、 $(n - fp - \lfloor f_L / 3 \rfloor) > n/3$ である。そこで次のようなアルゴリズムdfLEAD2でリーダー選択問題が解ける。(アルゴリズムdfLEAD2)

(1)起動プロセッサは”立候補者”となり、2個の

中継PEを作って自分の識別子を放送し、受信PEが

- a) "立候補者"でも"支持者"でもないか、
- b)より識別子の小さい"立候補者"または、その"支持者"ならば、それを自分の"支持者"にする。

このとき各NPEは一時に2つ以上の"立候補者"の"支持者"にならないようにする。このために、支持している"立候補者"を交えるには現在支持している"立候補者"に許可を求め、その返事がくることが条件になる。

(支持者の数+中継PEの数+1)が $n/3$ より多くなつた"立候補者"は"リーダー候補"になる。

"リーダー候補"となつたPEは、支持者に別のPEの支持者となる許可を与えない。これにより"リーダー候補"になるPEは高々2個になる。

(2)"リーダー候補"となつたPEは、deforestプロトコルをもちいて自分の識別子を放送し、 $n/2$ 個より多くのPEを(先の"支持者"と同じような)"投票者"とする事に成功するとリーダーになる。

(3)リーダーの識別子の放送をdeforestプロトコルで行なう。

(注意:ネットワークのサイズ n が小さいときには中継PEを作る時点で問題が解けている場合があるので、 n が4、5個のときには中継PEの個数を1個に、 n が3個のときには中継PEをつくらないようにする。)

dfLEAD2は2つの中継PEによる部分的な放送と、deforestプロトコルによる放送、そしてその2つの放送に対する返信によって成り立つ。「部分的な放送」に要するメッセージ数は各起動PEごとに $O(n)$ で、deforestによる放送は高々3回しか行なわれない。そしてこれらの放送に対する返信は、これらの放送に用いられるメッセージの定数倍になるので、アルゴリズム全体での通信計算量は、 $O(n+n \cdot k + fp \cdot f_L + f_L \cdot \log f_L)$ になる。
 $fp + f_L \geq n/4$ よりこれは

$O(n+k(fp+f_L)+fp \cdot f_L+f_L \cdot \log f_L)$ と表せる。

5. 3 アルゴリズムdfLEADの通信計算量

[定理5] dfLEAD1とdfLEAD2によって構成されたリーダー選択アルゴリズムdfLEADは、 $n-2 \geq fp+f_L$ かつ、 $fp < n/2$ の条件下で、CN-SDにおけるリーダー選択問題を通信計算量

$O(n+k(fp+f_L)+fp \cdot f_L+f_L \cdot \log f_L)$ で解く。□

一方、次の定理6により、この通信計算量は最適である。

[定理6] CN-SDにおけるリーダー選択問題の通信計算量の下界は、 $\Omega(n+k(fp+f_L)+fp \cdot f_L+f_L \cdot \log f_L)$ である。

(証明)各起動PEが $(fp+f_L)$ よりも少ないメッセージを送信するようなリーダー選択問題を解くアルゴリズムAが存在すると仮定する。起動PEが1個だけのときに最初に出すメッセージがすべて故障リンクまたはFPEに送信される場合があり、このときAはリーダー選択問題を解くことができない。従って、通信計算量の下界は $\Omega(k(fp+f_L))$ である。一方、リーダー選択問題はリーダーの識別子を全てのNPEが知って終了しなければならないので、放送問題と同じ下界 $\Omega(n+fp \cdot f_L+f_L \cdot \log f_L)$ が成り立つ。□

6. あとがき

リンクとPEの故障が混在する完全ネットワークにおける放送問題について考察し、方向感覚の有無が通信計算量のオーダーに影響を与えることを示し、その放送プロトコルdeforestを用いてリーダー選択問題を解決するアルゴリズムを設計した。本稿で述べた放送

プロトコルdeforestを利用して連結性関連問題など、さらに複雑な分散問題を解くことが今後の課題である。

謝辞

日頃から分散アルゴリズムについて有意義な御討論を頂く萩原兼一助教授、韓国広原大学朴政鎬博士に感謝致します。

参考文献

- (1)H.H.Abu-Amara:"Fault-tolerant distributed algorithm for election in complete networks",Tech. Rep. RUU-CS-87-10, Dept. of Electrical and Computer Eng., University of Illinois(Apr. 1987).
- (2)R.Bar-Yehuda and S.Kutten:"Fault tolerant distributed leader election with termination Detection in General Undirected Networks",Tech. Rep. #409,Department of Computer Sci.,Technion (Apr.1986, Revised August 1986).
- (3)R.Bar-Yehuda,S.kutten,Y.Wolfstahl and S.Zaks : "Making Distributed Spanning Tree Algorithms Fault-Resilient (Extended Abstract)",Proc.of STACS 87 pp.432-444 (Feb. 1987).
- (4)M.Y.Chan and F.Y.L.Chin : "Optimal Resilient Ring Election Algorithms", Proc. of 2nd International Workshop on Distributed Algorithms (LNCS 312),(Jul. 1987).
- (5)R.G.Gallager, P.A.Humblet and P.M.Spira : "A distributed algorithm for minimum-weight spanning trees", ACM TOPLAS,5,1, pp.66-77(Jan.1983).
- (6)O.Goldreich and L.Shrira:"Electing a Leader in a Ring with Link Failures", Acta Informatica,24, pp.79-91 (February 1987).
- (7)E.Korach,S.Moran and S.Zaks : "Tight lower and upper bounds for some distributed algorithms for a complete network of processors", Prc. of 3rd PODC, pp.199-207(Aug. 1984).
- (8)S.Kutten, Y.Wolfstahl and S.Zaks:"Optimal distributed t-resilient election in complete networks", Tech.Rep. #430,Dept. of Computer Sci., Technion(Aug. 1986).
- (9)M.C.Loui, T.A.Matsushita and D.B.West : "Election in a complete network with a sense of direction", Information Processing Letters,22,4, pp.185-187 (Apr. 1986).
- (10)増澤, 西川, 萩原, 都倉, 藤田:"方向感覚付き完全ネットワークにおけるリンク故障を考慮したリーダー選択問題", 信学技報, COMP88-98(1989.3).
- (11)D.Rotem, K.Korach and N.Santoro : "Analysis of a distributed algorithm for extrema finding in a ring", Journal of Parallel and Distributed Computing,4,6, pp.575-591 (Dec. 1987).
- (12)N.Santoro : "On the message complexity of distributed problems",International Journal of Computer and Information Sciences,13,3, pp.131-147 (Jul. 1984).
- (13)N.Santoro : "Sense of direction,topological awareness and communication complexity", SIGACT NEWS,16,2, pp.50-56 (1984).