

## リーダー選出問題における 時間最小アルゴリズムについて

京都大学工学部 大戸 豊 (Yutaka OHTO)  
京都大学工学部 茨木 俊秀 (Toshihide IBARAKI)

### あ ら ま し

分散システムにおいては、集中型のシステムにはなかった種々の問題が生じる。その一つとして、コントロールトークンの再生において、その作業をつかさどるプロセッサを定めたり、またファイル等の資源管理を指示するプロセッサを決定することがしばしば要求される。これが、リーダー選出問題である。本報告では、リーダーが選出されるまでの時間を減少させることに重点をおき、時間最小の分散アルゴリズムを提案する。

### 1. ま え が き

複数のプロセッサを通信リンクでつないで構成する"分散システム"は集中型システムに比べ、1)CPUやメモリ等の資源を共有できる、2)信頼性が向上する、3)負荷の分担が可能である、等の利点がある。ところが分散システムでは、集中型システムにはない新たな問題が発生してくる。分散システムには中央コントローラが存在しないため、例えば、システム内にコントロールトークンを流すことで各プロセッサの競合を制御する方式が考えられるが、そのコントロールトークンが消失した場合、新たにコントロールトークンをつくる"リーダー"と呼ばれる一つのプロセッサを、システム内から選び出す必要が生じる。あるいは、分散システム内のファイル等の資源管理をつかさどるプロセッサの決定が要求されることもある。これがリーダー選出問題である。

システムが非同期で、形態が一般グラフの場合、メッセージ数  $O(n \log n + m)$  [ $n, m$  はそれぞれシステム内のプロセッサ数とリンク数] の分散アルゴリズムが提案されている [1, 2, 3, 4]。この中で、文献 [1] のアルゴリズムは時間  $(n \log^* n)$  を持つ。

ところで、リーダー選出問題におけるこれらの研究では、リーダー決定までに必要なメッセージ数の最小化が中心に議論されてきた。それは、メッセージの伝送にかかる通信コストが計算コストに比べ高価であり、いかにメッセージを少なくするかということが実用上重要であるからである。しかし、最近光ファイバーなどの進歩により、通信リンクの容量は増加し、その結果、通信コストは低下の傾向にある。またLAN (ローカルエリアネットワーク) においては通信コストはほとんど無視できるという状況にあり、メッセージ数より、リーダー選出までの時間の減少がより重要な場合も考えられる。

そこで本報告においては、時間を減少させることに重点を置き、リーダー選出問題の時間最小アルゴリズムを提案する。そのために、まず、すべてのプロセッサにおいてシステムの形態を示すグラフを構築する問題を考え、一つの時間最小アルゴリズムを与え、このアルゴリズムがグラフ構築までに必要とするメッセージ数は、 $(2m^2 - mn + m)$ 以上、 $(2m^2 - mn + 2m)$ 以下であることを示す。グラフが構築できれば、その情報にもとずいてリーダーを選出できるから、このアルゴリズムはリーダー選出アルゴリズムでもある。本報告の後半では、グラフ構築問題の任意の時間最小アルゴリズムは、 $m(m+1)$ 個のメッセージを必要とすることを示す。したがって、上のアルゴリズムはメッセージ数のオーダーの意味で最良である。さらに、上のアルゴリズムがリーダー選出問題のアルゴリズムとしても時間最小であることを示す。

## 2. 諸定義と仮定

ここでは、準備段階として本論文で想定する分散システムに関する諸定義および仮定を行う。これらは、分散システムに対し、通常用いられるものである(例えば[6]参照)。本報告においては、分散システムにおけるプロセッサおよびそれらを結ぶリンクをそれぞれグラフ  $G=(V, E)$  で表現する。

$V$  : 節点(プロセッサ)の集合。

$E$  : 枝(リンク)の集合。

また、 $n=|V|$ 、 $m=|E|$ とする。以下、プロセッサを節点、リンクを枝と呼ぶ。

[定義1](分散システムにおける性質・仮定)

- 1) グラフ  $G$  の形態にはあらかじめ何の制限もおかない。但し、有限性および連結性は仮定する。またアルゴリズムの実行中に  $G$  は変化しない。
- 2) 中央コントローラは存在しない。
- 3) 節点はメッセージを送ることによってのみ通信できる。通信は、一般に非同期で行なわれる。またメッセージは誤りなしに送られる。
- 4) 同じ枝を通過してきたメッセージの順序は変わらない(FIFO)。また枝上で反対向きのメッセージがお互いに影響を及ぼすことはない。
- 5) 各々の節点は "id" と呼ばれる一意的な値を持っており(以下の議論では、簡単のため節点とそのidを同一視して扱う)、また何個の節点とつながっているかという次数  $deg$  を知っていて、それらを区別できる(すなわち接続枝(ポート)には番号  $1, 2, \dots, deg$  が付されている)。しかし、初期状態において各節点はこれら以外に何の情報も持たない。

- 6) 各々の節点は局所的メモリおよび、任意個のメッセージを保てるバッファを持つ。
- 7) 節点内での処理およびメッセージの発信受信に要する時間は、枝上の伝送の遅れと比較して無視できる。また、一つの枝上のメッセージの伝送時間は、それが発せられた時刻によって定まり、メッセージの内容にはよらない。
- 8) 2つ以上のメッセージが1つの節点に同時到着することはない(適当に順序をつけバッファに入れることができる)。□

ところで、 $n$ 個の節点の中から1つをリーダーとして選出するには、各節点の持つ一意的なidにもとずいて、その中で最大のidをもつ節点を選べばよい[5,6]。そのようなアルゴリズムを評価する際、アルゴリズムが始まってからリーダーが決まるまでに必要な時間およびメッセージ数が問題になる。そのためには、アルゴリズムの起動およびリーダーが決まる、ということの定義を明確にしておく必要がある。

[定義2](アルゴリズムの起動) 静止状態にある $n$ 個の節点のうち $p$ 個( $1 \leq p \leq n$ )がリーダー選出の必要性を知り、他にメッセージを送ることでリーダー選出アルゴリズムを起動する。リーダー選出の必要性を示すメッセージを受け取った節点は、そのことを他に知らせると同時に自己のリーダー選出アルゴリズムを起動する。□

[定義3](リーダー決定) リーダーが決まったとは次の2つをともに満足する時のことをいう。

- 1) 最大のidをもつ節点が、自分のidが最大であることを知る。
- 2) その他のすべての節点は、自分のidは最大ではなく、しかも最大のidが(つまり、リーダーが)みつかったことを知る。□

### 3. アルゴリズムとその考え方

節点 $k$  ( $k$ はid)において、 $ID_k$ 、 $deg_k(a)$ 、 $ADJ_k(a)$  (ただし、 $a \in ID_k$ )は、それぞれ、その時点で節点 $k$ がグラフ $G$ 中に存在することを知っている節点のidの集合、節点 $a$  ( $a \in ID_k$ )の次数、節点 $a$  ( $a \in ID_k$ )に隣接している(つまり、枝で結合されている)ことが判明している節点の集合である。各節点 $k$ が最初に受けるメッセージは $(i, p, d, \varepsilon)$ タイプであり、このメッセージは、1) 節点 $i$ の次数は $d$ である、2) 節点 $i$ はその $p$ 番目の枝によって節点 $k$ に隣接している、という情報を持つとともに、アルゴリズムの起動を知らせる役目も果たしている。 $\varepsilon$ はnull情報を意味する。もう一つのメッセージのタイプは、 $(i, p, d, j)$ であるが、これには、1) 節点 $i$ の次数が $d$ である、2) 節点 $i$ はその $p$ 番目の枝によって節点 $j$ に隣接している、という情報が含まれている。したがって、これらのメッセージを

受けると、各節点  $k$  ではデータ  $ID_k$ 、 $deg_k(a)$ 、 $ADJ_k(a)$ 、 $a \in ID_k$  を更新する。また、各節点  $k$  は、 $(i, p, d, j)$  (あるいは  $(i, p, d, \epsilon)$ ) のメッセージを受けたとき、ア) 枝  $(i, j)$  の情報 ( $(i, p, d, \epsilon)$  の場合は枝  $(i, k)$  の情報) をすでに知っている、イ) 自分は葉節点である、の場合以外は、そのメッセージをそのまま ( $(i, p, d, \epsilon)$  の場合は  $(i, p, d, k)$  に直して) 送られてきた以外のすべての方向へ送る。なお、ア) とイ) のどちらも成立しない場合の特別な場合として、 $i = k$ 、つまり自分の出したメッセージが一巡して戻ってきたという場合が考えられるが、このときには隣接する節点へはすでに  $k$  およびその次数の情報は送られているので、メッセージは、 $(i, \epsilon, \epsilon, j)$  でよい。特に、 $p$  番目の枝の相手は  $j$  であるから、 $j$  に対してこのメッセージを  $p$  番目の枝を通して送る必要はなく省略できる。

各節点  $k$  では、あるメッセージを受けて、その結果  $ID_k$ 、および  $a \in ID_k$  に対する  $deg_k(a)$ 、 $ADJ_k(a)$  からグラフの全体が構築できたとき、そのメッセージの他への発送をすませた後、アルゴリズムを停止する。このとき、 $ID_k$  にはグラフ内のすべての節点が入っているから、そのなかで最大のものをリーダーとすればよい。

リーダーが何であるかは、1) 構築されたグラフにもとづいて自分で決定する、2) 他節点から知らせてもらう、などの方法で、知ることができる。一番最初にリーダーがわかる節点は他節点から知らせてもらうことはできないから、自分で決定しなければならない。リーダーは最大の  $id$  をもつ節点であるから、リーダーを決定するには、グラフの全体が構築できればよい。節点  $k$  でのグラフの構築は、データ  $ID_k$ 、 $deg_k(a)$ 、 $ADJ_k(a)$  ( $a \in ID_k$ ) がすべて完成することと同値である。ここではまず全節点がグラフを構築するというグラフ構築問題を考え、その分散アルゴリズムを与える。

### グラフ構築アルゴリズム A

各節点  $k \in V$  はそれぞれ独立に以下のアルゴリズムを実行する。

- (1) 起動前:  $ID_k := \phi$ .
- (2) 起動:  $ID_k := \{k\}$ 、 $deg_k(k) :=$  (自分の節点の次数)、 $ADJ_k(k) := \phi$  とする。メッセージ  $(k, p, deg_k(k), \epsilon)$  を、自分に隣接する  $p$  番目の枝 ( $p = 1, 2, \dots, deg_k(k)$ ) を通して送る。 [ $\epsilon$  は null 情報を示す。]
- (3) メッセージ  $(i, p, d, j)$  を受け取った時 ( $p, d, j$  の一部が  $\epsilon$  である場合も含む):
  - Step 0 (起動):
    - (a) if  $ID_k = \phi$  [節点  $k$  はまだ起動していない。このとき受けとるメッセージは  $(i, p, d, \epsilon)$  タイプのみである]
    - (b) then  $ID_k := \{k\}$ 、 $ADJ_k(k) := \phi$ 、 $deg_k(k) :=$  (自分の節点の次数) とし、メッセージ  $(k, q, deg_k(k), \epsilon)$  を  $q$  番目の枝

( $q = 1, 2, \dots, \text{deg}_k(k)$ )を通して送る。

Step 1 (データの更新):

- (a) if  $j \neq \epsilon$ ,  $i \in \text{ID}_k$  and  $j \in \text{ADJ}_k(i)$  [メッセージ  
( $i, p, d, j$ ) or ( $j, p', d', i$ )をすでに受け取ったことがある]あるいは  $j = \epsilon$ ,  $i \in \text{ID}_k$  and  $k \in \text{ADJ}_k(i)$  [メッセージ ( $k, p'', d'', i$ )をすでに受け取っている]
- (b) then exit
- (c) else [枝 ( $i, j$ )の情報が新しく得られた]  
 $\text{ID}_k := \text{ID}_k \cup \{i\}$  および、 $d \neq \epsilon$  ならば  $\text{deg}_k(i) := d$  とする;  
 (i)  $j = \epsilon$  の場合  
 $\text{ADJ}_k(k) := \text{ADJ}_k(k) \cup \{i\};$   
 $\text{ADJ}_k(i) := \{k\};$   
 Step 2 へ。  
 (ii)  $j \neq \epsilon$  の場合  
 $\text{ADJ}_k(i) := \text{ADJ}_k(i) \cup \{j\};$   
 $\text{ADJ}_k(j) := \text{ADJ}_k(j) \cup \{i\};$   
 Step 2 へ。

Step 2 (メッセージの発送):

- (a) if  $\text{deg}_k(k) = 1$  [  $k$  は葉節点 ]
- (b) then Step 3 へ
- (c) else [  $k$  は葉節点ではない ]  
 if  $i = k$  [自分へ戻ってきたメッセージの処理]  
 then メッセージが送られてきた方向および、 $p$  番目の枝  
 を除いて、すべての枝にメッセージ  
 ( $i, \epsilon, \epsilon, j$ )を送る。  
 else [ $i \neq k$ ]  
 もし  $j = \epsilon$  ならば、 $j := k$  として、それ以外ならば  $j$  はそのままメッセージ ( $i, p, d, j$ ) を送られてきた以外の方向へ送る。Step 3 へ。

Step 3 (グラフ構築の完成テスト):

- (a) if  $|\text{ADJ}_k(a)| = \text{deg}_k(a)$  for all  $a \in \text{ID}_k$  [グラフの構築は完成]
- (b) then halt [以後、メッセージを受けても何の動作もしない]
- (c) else exit.  $\square$

なお、Step 3 の (b) は、アルゴリズムの停止を明示するために加えたもので、メッセージの発送に関しては影響を持たないことを注意しておく。この部分がなくても、グラフの構築が完成していれば、Step 1 (a) (b) によって、新しくメッセージを送り出すことはないからである。

#### 4. アルゴリズム A の正当性

枝の情報とは、その枝の端点を  $i, j$  としたとき、 $i$  と  $j$  が隣接しているという情報のことを言う。グラフを構築するには、全節点の次数と全枝の情報をすべて受け取らなくてはならない。

[補題1] アルゴリズム A によれば、すべての節点は有限時間内にすべての枝の情報を受け取ることができる。□

次に、節点の次数に関する情報を考える。各節点  $k$  はメッセージ  $(i, p, d, j)$  と  $(j, p', d', i)$  の一方のみを他へ発送する。その理由は、上の補題で述べたように、これらのメッセージは枝  $(i, j)$  について同等の情報を持つからである。しかし、前者には、 $i$  の次数のみが、また後者には  $j$  の次数のみの情報が含まれているので、これらのメッセージは、次数に関して同等の情報を持つわけではない。このため次の補題が必要である。

[補題2] アルゴリズム A によれば、節点  $k$  が他の節点  $k'$  から発せられた一つのメッセージ(最初  $(k', p', d', e)$  の形しており、その後  $(k', p', d', j)$  の形に変更される)を受け取ったとき、そのメッセージが通ってきた経路上のすべての節点について、それらの存在と次数の情報をすでに受けとっている。□

補題2において、 $k'$  の次数も当然  $k$  に伝わる。したがって、補題1と補題2をすべての  $k$  に適用すれば、次の定理が言える。

[定理1] 有限時間内にアルゴリズム A は停止し、すべての節点はグラフを構築できる。□

#### 5. 時間量の評価

すべての節点でグラフを構築するアルゴリズムの時間最小性を次のように定義する。

[定義4] (時間最小性) 同一の起動条件の下で、任意の節点  $s$  において、 $s$  が節点  $k$  の次数を知る時刻、および、枝  $e$  の情報を得る時刻が、任意の  $k, e$  について、他のどのアルゴリズムと比べても遅くないとき、その分散アルゴリズムは時間最小であるという。□

時間最小性を達成するアルゴリズムの一つとして、アルゴリズム  $A^*$  を考える。 $A^*$  は次のように動作する。

- ・起動はアルゴリズム A と同様に行なう。
- ・メッセージを受けとったら、自分のメモリ内に新しく加えて記憶す

る。その後直ちに、

- ・現在自分がもっている情報をすべて、隣接する節点到(メッセージが送られてきた節点も含めて)送る。

明らかにこれ以上早く情報を伝達することは原理的にありえない。任意の枝の情報、および節点の次数の情報について、各節点到最初に到着する時刻をアルゴリズム A と  $A^*$  で比較することで、次の定理が成り立つ。

[定理2] アルゴリズム A は時間最小である。□

システムは一般に非同期であり、上の結果は一般に成立するものであるが、グラフ構築に要する時間の簡単な評価を得るため、枝上の伝送にちょうど一単位時間かかる同期的な動作を仮定して時間量を求めてみる。G の 2 節点  $s$  と  $t$  の距離  $\text{dis}(s, t)$  を  $s$  と  $t$  を結ぶ最小枝数の経路の枝数と定義し、 $D$  をネットワーク G の直径、すなわち最も離れている 2 節点間の距離とする。

[補題3] システムが同期的に動作する場合、アルゴリズム A において、節点  $s$  から発せられたメッセージ  $(s, p, d, e)$  あるいは  $(s, p, d, s')$  が枝  $e = (i, j)$  を通るとすれば、 $\text{dis}(s, i) = \text{dis}(s, j) = D$  の場合には、 $(D + 1)$  単位時間に、それ以外の場合には、 $D$  単位時間以内に通っている。□

[定理3] システムが同期的に動作するとき、アルゴリズム A によってすべての節点がグラフを構築するまでに必要な時間は、一斉スタートの場合、 $L$  単位時間である。ただし  $L$  は次のように定義される。ある節点  $k$  に対し、 $\text{dis}(k, i) = D$ 、 $\text{dis}(k, j) = D$  をみたす枝  $(i, j) \in E$  が存在するとき、 $L = D + 1$ 。その他の場合、 $L = D$ 。□

## 6. メッセージ数の評価

次に、アルゴリズム A のメッセージ数を評価する。そのため、一つの枝  $(i, j)$  に関するメッセージ  $(i, p, d, j)$  あるいは  $(j, p', d', i)$  の伝送の様子に注目する。

[補題4] アルゴリズム A において、一つの枝  $(i, j)$  に関するメッセージ  $(i, p, d, j)$  あるいは  $(j, p', d', i)$  が G の同じ枝を同方向に 2 回以上流れることはない。□

[補題5] アルゴリズム A において、一つの枝  $(i, j)$  に関するメッセージは G の各枝に多くて 2 回(しかもそのときは向きが反対)流れるだけである。□

[補題6] アルゴリズムAにおいて、枝 $(i, j)$ と節点 $k$  ( $k \neq i, j$ )を考えると、 $k$ に接続する枝の中で、 $(i, j)$ に関するメッセージを $k$ から外へ送り出さないものがちょうど一本存在する。□

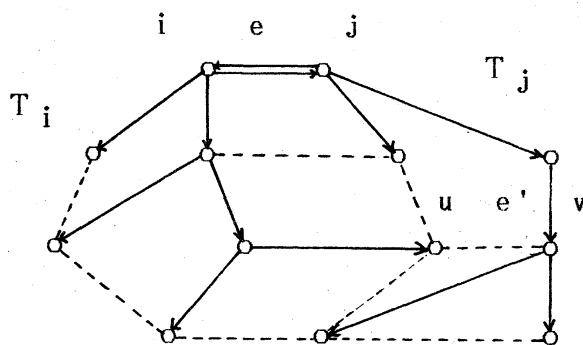


図1 枝 $e = (i, j)$ の情報を運ぶ木 $T_i$ と $T_j$

つまり、枝 $e = (i, j)$ に関するメッセージ $(i, p, d, j)$ と $(j, p', d', i)$ の伝搬形態に注目すると、 $(i, p, d, j)$ は節点 $j$ で作られ、 $j$ から $i$ 以外の節点へ送り出され、 $(j, p', d', i)$ は $i$ で作られ、 $i$ から $j$ 以外の節点へ送り出される。これらのメッセージが節点 $k$  ( $k \neq i, j$ )に最初に到着する経路は、それぞれ $j$ と $i$ を根とする2つの木 $T_j$ と $T_i$ によって表わされる(図1)。 $T_i$ と $T_j$ は互いに素でしか、両者と枝 $(i, j)$ を合わせると、 $G$ の全域木となっている。補題6によって $T_i$ と $T_j$ の枝には $(i, j)$ に関するメッセージがちょうど一回流れる。 $e' = (u, v)$ を $(i, j)$ でもなければ $T_i$ と $T_j$ のどちらにも属さない任意の枝とする(図1の点線の枝)。枝 $e'$ にはその両端 $u$ と $v$ から $(i, j)$ に関するメッセージが送られ途中で交差する。つまり、 $e'$ には $(i, j)$ に関するメッセージが2個流れる。最後に枝 $e = (i, j)$ には、 $(i, p, d, e)$ と $(j, p', d', e)$ の2個のメッセージが流れる(Step0(b))。これら全てをあわせると、メッセージ数は、

$$2(m - n + 2) + (n - 2) = 2m - n + 2$$

である。なお例外的な場合として、 $i$ から出たメッセージ $(i, p, d, e)$ に起動されて、Step0(b)によって $j$ で作られたメッセージ $(j, p', d', e)$ が $i$ に届く以前に、 $(i, p, d, j)$ が $i$ に到着する場合があります、このときはメッ



セージが最初に届いた経路が作る木は、 $j$ を根とする全域木一個である。このとき、メッセージ数は、

$$2(m - n + 1) + (n - 1) = 2m - n + 1$$

になる。アルゴリズムAのメッセージの送り方から、 $m$ 個の枝 $(i, j)$ すべてについて他の枝とは独立にメッセージが送られるので、結局次の定理が成り立つ。

[定理4] アルゴリズムAのメッセージ数は、 $2m^2 - mn + m$ 以上、 $2m^2 - mn + 2m$ 以下である。□

## 7. 時間最小のグラフ構築に必要なメッセージ数

以下の議論では、メッセージの単位として、枝一つおよびその両端点の情報(すなわち端点の次数、メッセージを送った枝番号)をとる。したがって、たとえば、一つの経路に関する情報は、このような単位に分解した形で送られるものとする。メッセージ数とは、このように分解されたメッセージの個数である。このとき、次の定理が成り立つ。

[定理5] 時間最小という条件のもとでは、全節点でグラフを構築する任意の分散アルゴリズムにおいて、メッセージ数 $m(m+1)$ は必要である。□

## 8. リーダー選出問題とグラフ構築問題

今までリーダー選出問題をグラフ構築問題におきかえた議論をしてきた。グラフが構築できれば、リーダーは選出できるから、アルゴリズムAを実行することで、有限時間内にリーダーを選出できる。しかし、リーダーは自分で決定しなくても、他から知らせてもらうことも可能であるから、全節点がグラフを構築しなくてもよい場合も考えられる。この場合、リーダー選出問題とグラフ構築問題では、メッセージ数、時間量に差がででくる可能性があるが、次の定理が成り立つ。

[定理6] リーダー選出問題はどの節点においてもグラフ構築問題よりも早く解けることはない。□

### 謝辞

日頃ご討論いただく研究室の諸氏に感謝いたします。なお、本研究は一部、文部省科学研究費によるものである。

文 献

- (1) Gafni, E.: "Improvements in the Time Complexity of Two Message-Optimal Election Algorithms", In Proc. 4th Ann. ACM Symp. Princ. Distr. Comp.:175-185(1985).
- (2) Afec, Y. and Gafni, E.: "Time and Message Bounds for Election in Synchronous and Asynchronous Complete Network", In Proc. 4th Ann. ACM Symp. Princ. Distr. Comp.:186-195(1985).
- (3) Gallager, R.G., Humblet, P.A. and Spira, P.M.: "A Distributed Algorithm for Minimum-Weight Spanning Trees", ACM TOPLAS 5:66-77(1983).
- (4) Korach, E., Moran, S. and Zaks, S.: "Tight Lower and Upper Bounds for Some Distributed Algorithms for a Complete Network of Processors", In Proc. 3rd Ann. ACM Symp. Princ. Distr. Comp.:199-207(1984).
- (5) LeLann, G.: "Distributed Systems - Towards a Formal Approach", in Information Processing 77:155-160(1977).
- (6) Liestman, A.L. and Peters, J.G.: "Distributed Algorithms", Technical Report TR86-10, School of Computing Science, Simon Fraser University, June(1986).