

分散型システムにおける相互排除のための トークンを用いたスキームについて

京都大学・工 室 章治郎(Shojiro Muro)

神戸商科大学 加藤 直樹(Naoki Katoh)

オレゴン州立大学 箕浦 敏美(Toshimi Minoura)

京都大学・工 長谷川利治(Toshiharu Hasegawa)

1. まえがき

相互排除 (mutual exclusion) 問題は、もともと、集中型
計算機システムにおいて、干渉しあう並行プロセスが、資源
(変数、ファイル等)を共用するために必要な排他制御を
実現する時間的な同期化規則として規定された [DIJK-65]。並
行プロセスは、臨界領域 (critical region, 以後, CRと略す)
とよばれる構造化された(プログラム)文の中でのみ共用資源
を参照し書きかえることができ [DIJK-65, HANS-73], CR文を
実行できる(実行がおわる)ことをCR内に入れる(CRから出
る)という。同じ共用資源を参照するCRが時間的な相互排除
を実現し、プロセス間の実行が滞りなく進むために、並行プ
ロセスは各共用資源に関する次の臨界領域条件を満たさなけ

ればならない。

(i) CRに入ろうとしたプロセスは、その要求が有限時間内に受け付けられる。

(ii) プロセスは、有限時間内でCRから出る。

(iii) 一時には高ターツのプロセスしかCRに入れない(相互排除の要請)。

分散型計算機システムは、1970年代におけるコンピュータの低価格化、通信技術の高度な発達という状況のもとで出現し、その制御法については現在も活発に研究が進められている。分散型システムにおいては、ある共有資源 ν (特に、そのコピーが複数個のサイトに置かれている場合)の臨界領域条件を満たすために、どのサイトからも参照できる共有メモリ(または、バッファ)を置き、(コピーも含めた)資源 ν に関する集中制御を行なうことは、システムの分散制御の思想からも好ましくなく、サイト間のメッセージ交換のみによる制御法が一般的に用いられている。今までに、論理時計を用いたLampertの方法[LAMP-78]をはじめとして、分散型システムに対して多くの相互排除アルゴリズムの提案およびその解析が行なわれているが[LELA-77, REED-79, PETE-79, RABI-80, RICA-81, SUZU-82], 特に, [SUZU-82]では、その論文で提案されたアルゴリズムも含めいくつかのアルゴリ

ズムについて、メッセージ交換により生ずるパスの構造、メッセージ数およびその種類を比較しながら議論してあり興味深い。

ところで、今までに提案された分散型システムに対する相互排除アルゴリズムでは、システム内に故障が発生した場合のアルゴリズムの継続的な実行に関してはあまり論じられなかった。本稿では、最初に、トークンの譲渡を基本とし、メッセージ交換によって実現される分散型システムの相互排除アルゴリズムを提案する。次に、トークンがシステムの故障により失われ得る可能性のある場合に、トークンを正常に再生するアルゴリズムについて述べる。さらに、システム内のサイトまたは通信路に故障が発生した場合における、提案したアルゴリズムの継続的実行について概略のみ述べる。

2. システムモデル

本稿で議論の対象とする分散型システムは、故障の生じていない場合、次の条件を満たすものとする。(1)システムは N 個のサイトからなり、サイト名を $1, 2, \dots, N$ とする。(2)各サイトは共通にアクセスできる記憶部分をもたないものとする。(3)任意の二つのサイト間には、通信経路(二つのサイト i, j を直接結ぶ通信回線を通信路(i, j)といい、いくつかの通信

路をつないだ通信回線を通信経路とよぶことにする)が存在し、各サイトはその通信経路を通じたメッセージ交換により、他のサイトと通信することができる。(4)発送されたメッセージが発送順に目的のサイトに着くという仮定はない。(5)2サイト間の通信遅延は有限であると仮定するが、遅延の上限は設定しない。(6)各サイトには、繰返しCRに入っている一つのプロセスがあり、各プロセスはCRへ入っている要求を一旦出したら、そのCR要求が認可され、CRへ入り、CRから出るまで次のCR要求を出すことができない。但し、各プロセスがCR内に入っている時間は有限であるとする。

以上の条件のもとで、次節で制御トークン(Control token)をサイト間で次々と譲渡(transfer)して、相互排除を実現するアルゴリズムCTTA(control token transfer algorithm)を記述する。

ここで、サイト間のトークンの譲渡に関して次の仮定をする。(7)トークンがあるサイトを出て、任意の他のサイトに行き着くまでの時間は有限である。

3. アルゴリズムCTTA

アルゴリズムCTTAでは、各共有資源 v に関してサイト間にトークンが1個存在し、トークンを有するサイトのプロセス

のみが資源 ν に関してCRに入ることができる。サイト i に n 回目のCR要求が発生したとき、他の各サイト j に送られるメッセージは $CRREQ(i, j, n)$ の形をしており、 n をサイト i の要求番号とよぶ。各サイト i には、長さ N の配列 LR_i があり、各 $j \neq i$ に対して $LR_i(j)$ はサイト i がその時点までに受け取ったサイト j からの $CRREQ$ メッセージの最大の要求番号が格納され、 $LR_i(j) < n$ なる $CRREQ(j, i, n)$ がサイト i に届くたびに $LR_i(j) \leftarrow n$ と更新される。 $LR_i(i)$ は、サイト i が最後に送った $CRREQ$ の要求番号を記憶する。トークンは長さ N の配列 HR をもち、 $HR(i)$ はその時点までにサイト i がCRに入った回数を示す。サイト i は $HR(i)$ の値を保持しているものとする(次節のトークン再生アルゴリズムで必要)。サイト i にCR要求が発生すると、 $LR_i(i) \leftarrow LR_i(i) + 1$ と更新し、メッセージ $CRREQ(i, j, LR_i(i))$ を他の各サイト j に送る。

CR要求を出していたサイト i にトークンが訪れた場合、サイト i はCRに入っている。CRから出るとき、 $HR(i) \leftarrow HR(i) + 1$ と更新し(この時点で、 $HR(i) = LR_i(i)$ が成立する)、

$$LR_i(j) > HR(j) \quad (3.1)$$

の条件を満たす j のうち N の剰余のもとで i の次に大きいサイト番号をもつサイトにトークンを譲渡する。

以上のようにしてトークンの譲渡が行なわれて、CR要求を

出していたサイトが次々とCRに入っていく。

(注意3.1) アルゴリズムCTTAは、次にCRに入るサイトの決定の仕方が異なることを除き、[SUZU-82]のAlg 1と本質的にはほぼ同じであり、[SUZU-82]のメッセージ"YOURTURN"をトークンと対応させることができる。[SUZU-82]では、条件式(3.1)を満たすサイトを調べるたびに、条件を満たすサイトを Q とよばれるqueueに記憶し、 Q に入ったサイト順にトークンが譲渡される。システムの満たす条件(4)を考えると、このような Q を用いたとしても、CR要求を出した時間的順序とCRに入るサイトの順序とは必ずしも一致せず、また、トークンが失なわれ Q の情報も失なわれた場合の損失の大きさを考え、アルゴリズムCTTAでは、トークンの譲渡に関して簡単な方法を採用した。 □

(注意3.2) 条件式(3.1)を満たす μ が存在しないとき、トークンはサイト i に停まっているものとする。サイト i も含めて、どれかのサイトがCR要求のメッセージを出せば、2節の条件(5)より、有限時間でサイト i にその要求が着き、トークンの譲渡が行なわれる。 □

共有資源 μ に関して、アルゴリズムCTTAではトークンを有するサイトのみがCRに入れれること、さらに、システムの満たす条件(5), (6), (7) および注意3.2より明らかに次の定理

が成立する。

定理3.1. アルゴリズムCTTAは、各共有資源に関して、臨界領域条件を満たす。 □

アルゴリズムCTTAでサイト i がCR要求を出してCRに入るまでに必要なメッセージ数は、[SUZU-82]のAlg 1 同様 N であり、[RICA-81]のアルゴリズムのメッセージ数の半数である。また、メッセージの種類は、 $LR_j(i)$, $HR(i)$ の i が可算無限個の値を取ることより、可算無限種必要である。

4. トークン再生アルゴリズム

2節で述べたシステムの満たすべき条件下では、トークンがシステム内で失なわれること(以後、トークンロスという)は考えられない。ところが各サイトや通信路に故障が生じるとトークンロスの可能性があり、トークンの再生が必要となる。各サイトが、例えばCR要求を出して一定時間待ってもトークンが譲渡されて来ないというような状況のもとで、トークンロスの疑いをもった場合、トークンの存在を調べ、もし、トークンが失なわれていたらトークンを1個再生するアルゴリズムREGENERATEを実行することができる。

各サイト i は、サイト i の年齢とよばれる非負整数 $CG(i)$ を、また、トークンは、トークン年齢とよばれる非負整数 TG

を保持し，これらの値は0に初期値設定されているものとする。トークン(再生した場合には，再生後のトークン)を保持しているサイト i に対して，常に $TG = CG(i)$ が成立するようにアルゴリズムが働く。

あるサイト i がトークンロスの疑いをもった場合，新しいトークン年齢の候補 g_i を次式で定める。

$$g_i = \min \{ l \mid l = nN + i \text{ (} n \text{ は非負整数) and } l > CG(i) \}.$$

$CG(i)$ を $CG(i) \leftarrow g_i$ と更新し，他の各サイト j ($j \neq i$) に $CG(j)$ を更新するためのメッセージ $UDCG(i, j, g_i)$ を送り，それらのサイトから返事が返ってくるのを待つ。もし，全てのサイト j から $ACK(j, i)$ が返ってきたらサイト i でトークンを再生し， TG を $TG \leftarrow g_i$ と更新する。もし，あるサイト j から $NACK(j, i)$ が返ってきたら，サイト i でのトークンの再生， TG の更新は行なぬ。

メッセージ $UDCG(i, j, g_i)$ を受け取った各サイト j では，もし， $g_i < CG(j)$ ならばサイト i に $NACK(j, i)$ を返す。 $g_i > CG(j)$ の場合，次の二通りに分かれる。もし，トークンがサイト j にあるならば， $CG(j) \leftarrow g_i$ ， $TG \leftarrow g_i$ と更新して， $NACK(j, i)$ をサイト i に返す。トークンがサイト j になければ， $CG(j) \leftarrow g_i$ と更新し， $ACK(j, i)$ を $HR(j)$ を付して(サイト i でトークンを再生するために必要)サイト i に返す。

以上のアルゴリズム REGENERATE において，トークンが複数個できるのを防ぐために，アルゴリズム CTТА に次の修正を加える必要がある。

サイト i にトークンが譲渡され，そのトークン年齢が TG であるとす。もし， $TG < CG(i)$ なら譲渡されてきたトークンを除去する。それ以外の場合，サイト i は CR へ入り，以後アルゴリズム CTТА に従ってトークンを譲渡していく。

各サイトで定義されるトークン年齢の候補 g_i の値が各サイトごとに異なり，しかも単調増加すること。サイト i がトークンを生成できるためには， g_i が全ての $CG(j)$ ($j \neq i$) より大きくなければならないこと。さらに，トークンが通信路上にあるときにどこかのサイトで新たなトークンが生成された場合，通信路のトークンがその後最初に訪れたサイトで，アルゴリズム CTТА の修正により，除去されることより次の定理を示すことができる。

定理 4.1. アルゴリズム REGENERATE の実行後必ずトークンが存在し，しかも 2 個以上トークンが作られることはない。

□

5. 故障が生じた場合のアルゴリズム CTТА の実行

最初に，本稿で対象とする故障の種類と故障時にシステムが満たすと仮定してよい条件について [GARC-82] を参考に述べる。

故障の種類：システム内の故障の種類として，サイトの故障と通信路の故障を考える。故障していないサイトを活性サイトといい，通信路の故障が原因である2サイト間を結ぶ全ての通信経路が正常なメッセージの伝送を行なえない場合，サイトの分割が起ったという。

故障時にシステムが満たす条件：サイト i または通信路 (i,j) の故障が発生した場合，サイト i または通信路 (i,j) の両端 i, j （通信路 (i,j) の故障中，サイト i, j は故障しないものとする）がそれぞれ責任をもち，有限時間内で他のサイトに故障の発生が知らされる（実現法は[MINO-82]に詳しい）。同時に故障中であるサイトの最大数は小さく，サイトの分割は高々二分割までである。サイト i に故障が起った場合，サイト i はメッセージを受信し，他のサイトに伝送することはできるが，メッセージの内容に応じて， $HR(i)$, $CG(i)$ および配列 LR_i の値は一切更新できない。また，これらの値は故障時に失なわれない（失なわれても，システム内のログ[GARC-82]を用いて再生可能である）。サイトおよび通信路の故障は有限時間で回復する。

サイトまたは通信路が故障した場合，故障発生時にトークンが失なわれなければ，活性サイト間では臨界領域条件が満たされるために，それぞれの故障の場合に対してアルゴリズム

△CTTA または REGENERATE を以下のように修正しなければならない。

サイトの故障時のアルゴリズムの修正：サイト i が故障前に CR 要求を出したために、トークンがサイト i に故障中に喪失されてきたら、そのトークンをサイト i と通信路で結ばれている任意のサイト j に喪失する。また、サイト i の故障中、 $LR_i(j)$ ($j=1, 2, \dots, N$) の値が更新されないため、各 $LR_i(j)$ の値が現在の $LR_j(j)$ の値に比べて非常に小さく、サイト i の故障回復直後にサイト i にトークンがある場合、式(3.1)を満たすサイトが一つもなくトークンの動きが止まる場合がある。そのため、CR 要求を出している各サイト j は、再度同じ $LR_j(j)$ の値をもって $CRREQ(i, j, LR_j(j))$ を出してもよいものとする。また、サイトの故障中、他の任意のサイト j からのメッセージ UDCG に対しては、無条件に NACK を返す。

通信路 (i, j) の故障によるサイト分割時のアルゴリズムの修正：サイト $i(j)$ は、サイト $i(j)$ と通信経路で結ばれているサイトの集合に、サイトの分割が起った旨のメッセージ PARTITION(i) (PARTITION(j)) を送る。そのメッセージを受け取った各サイトは、アルゴリズム CT TA は続行するが、アルゴリズム REGENERATE は実行せず、また、他の任意のサイトからのメッセージ UDCG を受け取っても無条件に

NACKを返す。サイトの分割により，CRREQが全てのサイトに届かない場合があり，サイトの故障の時と同様再要求を許す。

6. むすび

本稿では，システムの故障によりトークンが失われることも考慮した臨界領域条件を満たすアルゴリズムについて論じた。現在，故障時にシステムが満たす条件等を検討中である。

[謝辞] 末筆ながら，室，加藤が日頃御指導頂き，本稿に対してもコメントを頂いた京都大学 茨木俊秀 助教授に感謝する。

参考文献

- [DIJK-65] Dijkstra, E.W., "Solution of a problem in concurrent programming control", CACM, 8, 9 (1965), 569.
- [GARC-82] Garcia-Molina, H., "Reliability issues for fully replicated distributed databases", IEEE Computer (Sept. 1982), 34-42.
- [HANS-73] Hansen, P.B., "Operating system principles", Prentice-Hall, Inc. (1973).
- [LAMP-78] Lamport, L., "Time, clock, and the ordering of events in a distributed system", CACM, 21, 7 (1978), 558-565.
- [LELA-77] LeLann, G., "Distributed system - Toward a formal approach", IFIP 77 (1977), 155-160.
- [MINO-82] Minoura, T., "Resilient extended true-copy token scheme for a distributed database system", IEEE Tr. on SE (1982), 173-189.
- [PETE-79] Peterson, G.L., "Concurrency and complexity", TR-59, Department of Computer Science, The University of Rochester (1979).
- [RABI-80] Rabin, M.O., "N-process synchronization by $4 \cdot \log_2 N$ -valued shared variable", Proc. 21st IEEE Annual Symp. on FOCS (1980), 407-410.
- [RICA-81] Ricart, G. and Agrawala, A.K., "An optimal algorithm for mutual exclusion in computer networks", CACM, 24, 1 (1981), 9-17.
- [REED-79] Reed, D.P. and Kanodia, R.K., "Synchronization with event counts and sequences", CACM, 22, 2 (1979), 115-123.
- [SUZU-82] Suzuki, I and Kasami, T., "An optimality theory for mutual exclusion in computer networks", Proc. 3rd Inter'l Conf. on DCS (1982), 365-370