

共有メモリマルチプロセッサシステムにおける 同期時間最適な無待機時計合わせプロトコル

守屋 宣 井上美智子 増澤利光 藤原秀雄

奈良先端科学技術大学院大学 情報科学研究科
〒 630-01 奈良県生駒市高山町 8916-5
e-mail:{sen-m, kounoe, masuzawa, fujiwara}@is.aist-nara.ac.jp

あらまし 共有メモリマルチプロセッサシステム、特に、システム内の全プロセッサが大域パルスを共有するフェーズ内システムにおける無待機時計合わせプロトコルに関して考察する。無待機時計合わせプロトコルとは、ある特定時間である同期時間以上正常に動作し続けているすべてのプロセッサの局所時計が同期することを保証するプロトコルで、プロセッサが任意の時間動作を停止し、動作停止に気づかず動作を再開する、居眠り故障に耐性がある。フェーズ内システムにおける無待機時計合わせプロトコルとしては、これまで同期時間 $O(n^2)$ のプロトコルが提案されていた (n はプロセッサ数)。本稿では、同期時間 $O(n)$ のプロトコルを提案し、また、この同期時間がオーダー的に最適であることを示す。

キーワード 共有メモリマルチプロセッサシステム、時計合わせプロトコル、無待機性、同期時間

Optimal Wait-free Clock Synchronization Protocol on a Shared-memory Multi-processor System

Sen MORIYA Michiko INOUE Toshimitsu MASUZAWA Hideo FUJIWARA

Graduate School of Information Science
Nara Institute of Science and Technology
8916-5, Takayamacho, Ikoma, Nara 630-01
e-mail:{sen-m, kounoe, masuzawa, fujiwara}@is.aist-nara.ac.jp

Abstract We consider wait-free clock synchronization protocols on a shared-memory multi-processor system, especially, on an in-phase system in which all processors share a common clock pulse. A wait-free clock synchronization protocol guarantees that, for a fixed k , local clocks of processors which have been working correctly for at least k time are synchronized. Such k is called synchronization time. The best previous protocol has synchronization time $O(n^2)$, where n is the number of processors. In this paper, we present a wait-free synchronization protocol with synchronization time $O(n)$. We also show that this synchronization protocol is asymptotically optimal.

key words shared-memory multi-processor system, clock synchronization, wait-freedom, synchronization time

1 まえがき

複数のプロセッサがそれぞれの局所時計を管理している分散システムを考える。このような分散システムにおいて、故障が存在するにも関わらず局所時計を同期させ続ける問題である時計合わせ問題は多くの応用問題にとって重要であり、またそれ自体も興味深い問題である。時計合わせ問題は、これまで、メッセージパッシングシステム、共有メモリマルチプロセッサシステム等、種々の分散システムにおいて研究されている。

メッセージパッシングシステムにおける時計合わせ問題は、各プロセッサがそれぞれ進捗率に偏りのある物理時計を持ち、かつ、メッセージの伝送遅延が有界であるという仮定のもとで、プロセッサのビザンティン故障、認証ビザンティン故障に耐性を持つプロトコルに関する多くの研究がなされた[1, 2, 3, 4, 5]。これらの研究で提案された多くのプロトコルは、正常なプロセッサの局所時計を同期させるもので、一度故障したプロセッサが復旧後に他の正常なプロセッサと局所時計を同期させることは、保証されていない。

また、Gouda ら[6]、Arora ら[7]は、分散システムの任意の状況から実行を開始しても、各プロセッサの時計を同期させることができる、自己安定時計合わせプロトコルを提案した。これらのプロトコルは、局所時計の破壊などの一時故障が起きても、その後、十分長い時間故障が起きなければ、時計の値を同期させることができる。

共有メモリマルチプロセッサシステムにおける時計合わせ問題は、共有メモリマルチプロセッサシステムにおけるプロセッサの動作と故障パターンを考慮して、Dolev と Welch によって提案された[8]。彼らは、プロセッサの故障として、任意時間動作を停止し、動作停止に気づかずして動作を再開するという居眠り故障を対象とし、任意個のプロセッサの居眠り故障に耐性のある時計合わせプロトコルを提案した。提案されたプロトコルは、ある固定された値 k に対し、以下の 2 つの条件を保証している。

- 少なくとも k 時間連続して動作したプロセッサは、それ以降正常に動作し続ける限り、局所時計の調整を行わない。
- 少なくとも k 時間連続して動作した 2 つのプロセッサの局所時計の値は一致する。

このプロトコルは、各プロセッサは少なくとも k 時間連続して動作すれば、他のプロセッサの動作に関わら

ず、局所時計を同期させることができるという意味で、無待機プロトコルと呼ばれる。また、上記の k を無待機時計合わせプロトコルの同期時間と呼ぶ。

本稿では、共有メモリマルチプロセッサシステム、特に、全プロセッサが共通の大域パルスを共有する、フェーズ内システムにおける無待機時計合わせプロトコルを考察する。フェーズ内システムにおける無待機時計合わせプロトコルとして、Dolev と Welch は、同期時間 $O(n^3)$ の自己安定無待機時計合わせプロトコルおよび、同期時間 $O(n^2)$ の無待機時計合わせプロトコルを提案した[8]。また、Papatriantafilou と Tsigas は、同期時間 $O(n^2)$ の自己安定無待機時計合わせプロトコルを提案した[9]。本稿では、同期時間 $O(n)$ の無待機時計合わせプロトコルの提案をする。また、提案するプロトコルの同期時間がオーダー的に最適であることを示す。

以下、2 節ではシステムモデルに関する定義を行う。3 節では、無待機時計合わせプロトコルの定義を行う。4 節では、同期時間 $O(n)$ の無待機時計合わせプロトコルを示し、プロトコルの証明の概略を示す。5 節では、無待機時計合わせプロトコルの同期時間の下界が $\Omega(n)$ であることを示す。

2 モデル

n 個のプロセッサがいくつかのレジスタを共有している共有メモリマルチプロセッサシステムを考える。プロセッサはレジスタを介してのみ通信を行うことができる。本稿では、各レジスタは $1WnR$ レジスタ (single-writer n -reader register) であるとし、各レジスタ R に対し、1 プロセッサのみが書き込みができる、全プロセッサが読み出しをできる。レジスタ R に書き込みができるプロセッサを R の所有者と呼ぶ。

各プロセッサは状態機械としてモデル化される。状態 s のプロセッサ P は、以下の操作を行い、次の状態 s' に遷移する。

- (1) 状態 s により決定するレジスタ R を読み込む。
- (2) 状態 s とレジスタ R の値を基に、 P が所有するレジスタを更新し、状態 s' に遷移する。

これらの操作を 1 ステップと定義する。

本稿では、フェーズ内システムと呼ばれる同期式マルチプロセッサシステムを扱う。フェーズ内システムとは、全プロセッサが共通の大域パルスを共有するシステムであり、全プロセッサは 1 ステップずつ同期し

て動作する。ただし、本稿では居眠り故障を考慮し、各パルスで全プロセッサが動作するとは限らないとする。このようなシステムを以下のようにモデル化する。

システムの大域状況は、全プロセッサの状態と全レジスタの値の組で表すことができる。以降、このシステムの大域状況のことを、単に状況と言う。システムの実行は、状況と、各パルスで動作したプロセッサの集合の無限または有限交互列 $E = c_0\pi_1c_1\dots$ で表すことができる。ここで、各 c_t は状況を表し、 c_0 は全プロセッサの初期状態、全レジスタの初期値から成る初期状況を表す。 π_t は t 番目のパルス発生時に動作したプロセッサの集合を表す。また、実行 E が有限であるときは、 E は状況で終る。この実行 E は、各 t に対し、 π_t に属する各プロセッサが c_{t-1} を基に 1 ステップを行い、状況が c_t になったことを意味する。以下では、パルスをこのプロセッサ集合として扱う。また、 t 番目のパルスが発生した（大域）時刻を時刻 t と呼ぶ。プロセッサ P について、 $P \notin \pi_t$ のとき、 P は時刻 t で居眠りしたと言う。時刻 t で居眠りしたプロセッサ P は c_{t-1} を基にしたステップは行わないため、 c_{t-1} と c_t における P の状態と P が所有するレジスタの値は変わらない。

以下では、システムの n 個のプロセッサを P_1, P_2, \dots, P_n と表し、各 P_i に対し i を識別子と言う。また、簡単化のため、本稿では各プロセッサ P_i が所有するレジスタは 1 個だけとし、そのレジスタを R_i と表す。 P_i は R_i をいくつかのフィールドに分割して使用することができる。以下では、レジスタ R のフィールド f を $R.f$ 、プロセッサ P の変数 x を $P.x$ と表す。実行 $E = c_0\pi_1c_1\dots$ に対し、状況 π_t でのフィールド $R.f$ 、変数 $P.x$ の値を、それぞれ $R.f(\pi_t), P.x(\pi_t)$ で表す。ただし、文脈から明らかなときは、単に $R.f(t), P.x(t)$ と表すこともある。

3 無待機時計合わせプロトコル

本節では、無待機時計合わせプロトコルを定義する。各プロセッサ P_i の所有するレジスタ R_i はフィールド $clock$ を持つ。 $R_i.clock$ の値は、プロセッサ P_i の局所時計の値を保持する。実行 $E = c_0\pi_1c_1\dots$ に対して、プロセッサ P が時刻 t まで連続して動作したパルス数を $work(P, t)$ と表す。すなわち、 $work(P, t) = \max\{l | P \in \bigcap_{t-(l-1) \leq t' \leq t} \pi_{t'}\}$ とする。ただし、 $P \notin \pi_t$ のときは、 $work(P, t) = 0$ と定義する。

任意の t に対し、時刻 t まで連続して動作したパル

ス数がある特定パルス以上であるような任意のプロセッサ P の局所時計の値が一致し、以降 P が動作し続ける限り局所時計の値が 1 パルス毎に 1 ずつ増えるようにするプロトコルを無待機時計合わせプロトコルと言う。無待機時計合わせプロトコルは、次のように定義される。

定義 1 ある正整数 k に対し、システムの任意の実行が次の 2 つの条件を満たすようなプロトコルを、無待機時計合わせプロトコルと言う。

Adjustment: 任意の $t \geq 1$ 、任意のプロセッサ P_i に対し、 $work(P_i, t) > k$ ならば、

$$R_i.clock(t) = R_i.clock(t-1) + 1$$

が成り立つ。

Agreement: 任意の $t \geq 1$ 、任意のプロセッサ P_i, P_j に対し、 $work(P_i, t) \geq k, work(P_j, t) \geq k$ ならば、

$$R_i.clock(t) = R_j.clock(t)$$

が成り立つ

また、このときの正整数 k を同期時間と言う。□

4 プロトコル

本節で、同期時間 $O(n)$ の無待機時計合わせプロトコルを提案する。プロトコルを図 1 に示す。

4.1 プロトコルの概略

本プロトコルでは、各プロセッサ P_i は 2 つのモード、調整中モード、調整済モードを持つ。調整中モードでは時計の調整を行い、調整済モードでは 1 ステップ毎に局所時計 $R_i.clock$ の値を 1 ずつ増やす。また、モードに関わらず、各ステップで、 P_i 自身あるいは他のプロセッサが居眠りしたかどうかのチェックを行う。

変数

各プロセッサ P_i は、 P_i の変数および R_i のフィールドとして、以下の値を保持する。

- $clock$ ：時計の値を表す。（初期値は 0）
- $clock_gen$ ：時計の世代番号を表す。局所再設定をする度に 1 増やされる。（初期値は 1）

- *work_count*: 現在の時計の世代番号の調整作業を始めてから、動作したステップ数を表す。(初期値は 1)
- *count*: 初期状況からのステップ数を表す。(初期値は 1)
- *adjusted*: 時計の調整を終えたかどうかを表す論理型変数で、時計調整中のときは 0、時計調整が終したら 1 である。(初期値は 0)
- *invalid[1..n]*: *invalid[j]* は R_j の時計の無効な世代番号のうち、 P_i が知っている最大値を表す。(初期値は 0)

また、上記以外に P_i の変数として、以下の値を保持する。

- *last[1..n]:last[j]* は 4 つのフィールド *my_count*, *count*, *clock_gen*, *work_count* を持ち、それぞれ前回 R_j の読み出しをしたときの $P_i.count$, $R_j.count$, $R_j.clock_gen$, $R_j.work_count$ の値を表す。(初期値は (1,1,1,1))
- *stay*: 時計調整を終えるまでの足踏み(後述)の回数を表す。(初期値は 0)
- *j*: 現在のパルスで読み出したレジスタが R_j であることを表す。(初期値は 1)

動作

調整中モードでは、以下のように局所時計を調整する。以下の場合、 P_i は調整作業を開始する。

- 初期状態である場合。
- P_i 自身の居眠りを検知して、または、調整時間が長くなり過ぎて、*clock*, *work_count* 等のいくつかの変数を再設定した場合。

後者の場合のいくつかの変数の初期化を局所再設定と呼ぶ。

初期状態、または局所再設定を行った後、プロセッサはまず $R_1, R_2, \dots, R_n, R_1, R_2, \dots, R_n$ の順に、プロセッサの識別子の順で 2 周、各ステップで 1 レジスタずつの読み出しを行う。レジスタ R_j の読み出しをしたとき、 $R_j.work_count$ の値は *last[j].work_count* に格納して記憶する。2 周目の読み出しを終えたとき、各レジスタから読み出した *work_count* の値を基に、各プロセッサを *work_count* の降順にソートす

る。ソートの結果を Q とする。3 周目は、自分より *work_count* が大きいレジスタに対し、系列 Q の順にレジスタの読み出しを行いながら、時計の値を調整する。このとき、 R_j の状態により、次のように操作する。

- 以下の場合、 $R_j.clock$ の値は無効とし、自分の *clock* の値を 1 増やす。
 - (1) 前回の R_j の読み出し以降に P_j が局所再設定を行っている場合
(P_j の時計の世代番号で判断する。)
 - (2) 前回の R_j の読み出し以降の P_j の居眠りを検知した場合
(居眠りのチェック法は後述する。)

- 上記以外の場合、次のようにする。ソートの結果 Q において、 P_i より先にあるすべてのプロセッサ P に対し、 P の *clock* を無効と判定する、あるいは P の *clock* に P_i の *clock* を合わせる操作を行ったら、 P_i は調整済とする。

- (1) P_j が調整済の場合、 P_i の *clock* の値を $R_j.clock$ の値に合わせる。ただし、この調整作業中に少なくとも 1 回他のプロセッサの局所時計に P_i の時計を合わせたことがある場合、以下のときは $R_j.clock$ に P_i の *clock* を合わせることはない。

- $P_i.clock > R_j.clock$ である場合は、 $R_j.clock$ の値は無効とし、自分の *clock* を 1 増やす。
- $P_i.clock < R_j.clock$ である場合は、 P_i 自身の居眠りを検知したことになるので、局所再設定を行う。

- (2) P_j が調整中の場合は、次ステップでも R_j を読み出すことにする。この操作を足踏みと呼ぶ。この調整作業における足踏み回数(変数 *stay*)が n になった場合は、 P_i は局所再設定を行う。

調整済モードでは、まず調整中モードの 3 周目で読み出しをしなかったレジスタを系列 Q の順に読み出しを行い、そのあとは系列 Q の順に読み出しを続け、各ステップでは、*clock* の値を 1 増やす。

居眠りのチェックは以下のようにして行う。プロセッサ P_i は、あるレジスタ R_j の読み出しをしてから再び R_j の読み出しをするまで、 $P_i.count, R_j.count, R_j.clock_gen$ の値をそれぞれ $last[j]$ のフィールド $my_count, count, clock_gen$ に格納して記憶する。レジスタ R_j の読み出しをしたとき、前回の R_j の読み出し時からの $R_j.count, P_i.count$ の増分の差を $dif_DeltaCount$ とする。このとき、 R_j の前回の読み出しから今回の読み出しの間で、 $dif_DeltaCount > 0$ ならば P_j が、 $dif_DeltaCount < 0$ ならば P_i が、居眠りしたことわかる。 $dif_DeltaCount < 0$ のときは、 P_i は局所再設定をし、調整作業を始めからやり直す。 $dif_DeltaCount > 0$ のときは、 $R_i.invalid[j]$ に $R_j.clock_gen$ の値を書き込み、 P_j が居眠りした世代を知らせる。また、逆に $R_j.invalid[i] = clock_gen$ の場合は、 P_i が世代番号 $clock_gen$ での居眠りを検知したことを意味し、 P_i は局所再設定を行う。

P_i は、各ステップの最後に、変数 $P_i.clock, P_i.clock_gen, P_i.work_count, P_i.count, P_i.adjusted, P_i.invalid[j]$ の値を、それぞれに対応する R_i のフィールドに書き込む。

4.2 プロトコルの正当性

図 1 のプロトコルが、同期時間 $O(n)$ の無待機時計合わせプロトコルであることを示す。ただし、紙面の都合上、いくつかの補題に対しては証明のアイデアのみを示す。

任意の実行 $E = c_0\pi_1c_1\cdots$ に対し、条件 Adjustment と Agreement が成り立つことを示す。2つの系列 s_1, s_2 に対し、 s_2 が s_1 からいくつかの要素を削除して得られる系列であるとき、 s_2 は s_1 の順序を保存すると言ふ。

プロセッサ P_i が時刻 t_s で手続き $sort$ を実行する。 P_i が、ソートして得た系列 $P_i.Q(t_s) = (P_{q_1}, P_{q_2}, \dots, P_{q_n})$ に対し、その接頭部系列 $(P_{q_1}, P_{q_2}, \dots, P_l)$ を $P_i.Q(t_s)^l$ と記す。また、時刻 t_s 以降、 P_i が居眠りおよび局所再設定を行なわずに調整を終えた（変数 $adjusted$ の値を 1 にした）場合に、調整を終えるまでに局所時計の値を合わせたプロセッサの集合を $A(P_i, t_s)$ と表す。 $A(P_i, t_s)$ に属するプロセッサ P_j は、 P_i が局所時計の値を合わせたときに $P_j.adjust = 1$ があるので、それ以前に最後に調整を終えた時刻 t_j^1 、最後に手続き $sort$ を行なった時刻 t_j^2 が存在する。 $A(P_i, t_s)$ に属するプロセッサ P_j のうち、 t_s より後で P_i が P_j の局所時計の値に合わせる以前に調整を終える

プロセッサ、すなわち、 $t_s < t_j^1$ であるプロセッサの集合を $A'(P_i, t_s)$ とする。このとき、 $work(P_i, t_s) > 5n$ であれば、任意の $P_j \in A(P_i, t_s)$ に対し、 $P_i.Q(t_s)^j$ は $P_j.Q(t_j^2)^j$ の順序を保存すること、また、任意の $P_j, P_{j'} \in A(P_i, t_s)$ に対し、 $P_i.Q(t_s)$ において P_j が $P_{j'}$ に先行すれば、 $P_{j'}.Q(t_{j'}^2)^j$ は $P_j.Q(t_j^2)^j$ の順序を保存することが保証される。このことから、 P_i が時刻 t_s 以降、最初に時計の調整を終える、または、局所再設定を行なうまでに行なう足踏みは高々 $n - 1$ であることが導ける。

$work(P_i, t_s + 2n) > 7n$ の場合を考える。このとき、 P_i は区間 $[t_s, t_s + 2n]$ では、自身の居眠りを検知して局所再設定を行なうこともない。このことから、 P_i は区間 $[t_s + 1, t_s + 2n]$ 中に調整を終える。すなわち、以下の補題が成り立つ。

補題 1 プロセッサ P_i が時刻 t_s で手続き $sort$ を実行する。 $work(P_i, t_s + 2n) > 7n$ ならば、

$$P_i.adjusted(t_s + 2n) = 1$$

である。 \square

補題 2 プロセッサ P_i が時刻 t_r で局所再設定する。 $work(P_i, t_r + 10n) > 10n$ ならば、

$$P_i.adjusted(t_r + 10n) = 1$$

である。

(証明) $t_r + 5n$ までに調整が終らない場合、すなわち、 $P_i.adjusted(t_r + 5n) = 0$ のときを考える。このとき、区間 $[t+4n, t+5n]$ 中のある t'_r で $stay = n$ となり、 P_i は局所再設定を行う。 P_i は $T'_r + 3n$ で再び手続き $sort$ を行なうが、このとき $work(P_i, (t'_r + 3n) + 2n) > 7n$ となるので、 $T'_r + 3n + 2n < T_r + 10n$ までに調整を終える。 \square

P_i が、時刻 t で R_j を読み出したときに、自分の居眠りの検知による局所再設定を行うのは、次の場合である。

- $dif_DeltaCount < 0$ のとき
- $R_j.invalid[i](t-1) \geq P_i.clock_gen(t-1)$ のとき
- $0 < P_i.clock(t-1) < R_j.clock(t-1)$ かつ
 $P_i.invalid[j](t) < R_j.clock_gen(t-1)$ のとき

それぞれの場合に, $\text{work}(P_i, t) < 7n$ が成り立つ. よって, 以下の補題が成り立つ.

補題 3 P_i が時刻 t で局所再設定するならば,

$$\text{work}(P_i, t) < 7n$$

が成り立つ. \square

補題 4 任意の $t \geq 1$, 任意のプロセッサ P_i に対し, $\text{work}(P_i, t) > 17n$ ならば,

$$P_i.\text{clock}(t) = P_i.\text{clock}(t - 1) + 1$$

が成り立つ.

(証明) P_i が t 以前で最後に居眠りした時刻を t_ν とする.

まず, P_i が t_ν 以降に居眠りの検知による局所再設定をしない場合を考える. このとき, $t_\nu + 5n$ までに調整を終えないならば, 補題 2 と同様に, 遅くとも $t_\nu + 10n$ までに調整を終える.

次に, P_i が t_ν 以降に居眠りの検知による局所再設定をする場合を考える. このとき, 補題 2 より, 局所再設定を行うのは遅くとも $t_\nu + 7n$ であり, それ以後は居眠り検知による局所再設定はしない. よって, 補題 3 より, P_i は遅くとも $t + 17n$ までに調整を終える.

P_i は, 調整を終えたあとは局所再設定しないので, プロトコルより, 各ステップで clock の値を 1 ずつ増やす. よって, $\text{work}(P_i, t) > 17n$ のときは

$$P_i.\text{clock}(t) = P_i.\text{clock}(t - 1) + 1$$

が成り立つ. \square

次に条件 Agreement が成り立つことを示す.

補題 5 任意の $t \geq 1$, 任意のプロセッサ P_i, P_j に対し, $\text{work}(P_i, t) \geq 17n, \text{work}(P_j, t) \geq 17n$ ならば,

$$P_i.\text{clock}(t) = P_j.\text{clock}(t)$$

が成り立つ.

(証明) 補題 2, 3 より, $P_j.\text{adjusted}(t) = 1$ が成り立つ. よって, P_i, P_j とともに, t 以前に調整作業を終えたことがある. P_i, P_j が t 以前に終えた最後の調整作業で局所時計を合わせたプロセッサの集合を, それぞれ, $\hat{A}(P_i, t), \hat{A}(P_j, t)$ とする. このとき, $P_j \in \hat{A}(P_i, t)$, または $P_i \in \hat{A}(P_j, t)$ が成り立つ (証明略).

一般性を失うことなく, P_i を,

(1) $P_j \in \hat{A}(P_i, t)$ かつ $P_i \notin \hat{A}(P_j, t)$ であるプロセッサ,

または,

(2) $P_i \in \hat{A}(P_j, t)$ かつ $P_j \in \hat{A}(P_i, t)$ なら, t 以前に終えた最後の調整作業で, P_j が P_i の局所時計に合わせた時刻 $t_{j,i}$ と P_i が P_j の局所時計に合わせた時刻 $t_{i,j}$ に対し, $t_{j,i} \leq t_{i,j}$ が成り立つプロセッサ

であるとする. このとき,

$$P_i.\text{clock}(t_{i,j}) = P_j.\text{clock}(t_{i,j} - 1) + 1$$

が成り立つ.

時刻 $t - 17n + 2$ 以降に P_i が最初に R_j を読み出す時刻と $t_{i,j}$ の最大値を $t'_{i,j}$ とする. P_i, P_j は $t'_{i,j}$ 以降, 局所再設定をすることはないので, $t_{i,j} < t'_{i,j} - 1$ のとき, 区間 $[t_{i,j}, t'_{i,j} - 1]$ で P_i が居眠りする回数を nap_i , P_j が居眠りする回数を nap_j とすると, $\text{nap}_i = \text{nap}_j$ が成り立つ. P_i, P_j は区間 $[t_{i,j}, t'_{i,j} - 1]$ で、それぞれ、各ステップで時計の値を 1 ずつ増やすだけであるので,

$$P_i.\text{clock}(t'_{i,j}) = P_j.\text{clock}(t'_{i,j} - 1) + 1$$

が成り立つ. また, P_i, P_j ともに区間 $[t'_{i,j} - 1, t]$ では居眠りをしないので,

$$P_i.\text{clock}(t) = P_j.\text{clock}(t)$$

が成り立つ. \square

補題 4 と補題 5 より, 次の定理が成り立つ.

定理 6 図 1 のプロトコルは, 同期時間 $O(n)$ の無待機時計合わせプロトコルである. \square

5 同期時間の下界

前節で, 同期時間が $O(n)$ の無待機時計合わせプロトコルを提案した. 本節では, 無待機時計合わせプロトコルの同期時間の下界が $\Omega(n)$ であることを示す.

定理 7 同期時間が $n - 2$ 以下の無待機時計合わせプロトコルは存在しない。

(証明) 同期時間が $k \leq n - 2$ の時計合わせプロトコル A が存在すると仮定し, 矛盾を導く. 以下の 2 つの条件を満たすプロトコル A の 2 つの実行 $E = c_0\pi_1c_1 \dots, E' = c'_0\pi'_1c'_1 \dots$ を考える.

- E の c_t までの接頭部は、 E' の c'_t までの接頭部と一致する。
- 任意のプロセッサ P に対し、 $\text{work}(P, t) \geq k$ が成り立つ。

プロトコル A は同期時間 k のプロトコルであるので、全プロセッサ P_i に対し、 $R_i.\text{clock}(t)$ が一致する。この値を $\text{clock}(t)$ とする。

実行 E において、 $\pi_{t+1} = \pi_{t+2} = \dots = \pi_{t+n-1} = \{P_1\}$ とする。このとき、条件 Adjustment より

$$R_1.\text{clock}(c_{t+n-2}) = \text{clock}(c_t) + n - 2 \quad (1)$$

が成り立つ。実行 E では、 P_1 は時刻 $t+1$ から $t+n-2$ までに $n-2$ ステップの動作を行うので、その間に読み出しているレジスタが少なくとも 1 つ存在する。そのレジスタを R_i とする。

実行 E' において、 $\pi'_{t+1} = \{P_i\}$, $\pi'_{t+2} = \pi'_{t+3} = \dots = \pi'_{t+n-1} = \{P_1, P_i\}$ とする。このとき、条件 Adjustment より

$$R_i.\text{clock}(c'_{t+n-1}) = \text{clock}(c'_t) + n - 1 \quad (2)$$

が成り立つ。また、 P_1, P_i 以外のプロセッサは、 $\pi_{t+1}, \dots, \pi_{t+n-1}, \pi'_{t+1}, \dots, \pi'_{t+n-1}$ のどのパルスにも属していない。よって、 R_1, R_i 以外のレジスタの値は実行 E では c_t から c_{t+n-1} まで、実行 E' では c'_t から c'_{t+n-1} まで更新されない。故に、 P_1 は実行 E' での時刻 $t+2$ から $t+n-1$ では、 E の時刻 $t+1$ から $t+n-2$ までと全く同じステップを行う。従って、式(1)より実行 E' において、

$$R_1.\text{clock}(c'_{t+n-1}) = \text{clock}(c'_t) + n - 2 \quad (3)$$

が成り立つ。ところが、実行 E' では $\text{work}(P_1, t+n-1) = n-2 \geq k$, $\text{work}(P_i, t+n-1) = t+n-1 \geq k$ なので、条件 Agreement より

$$R_1.\text{clock}(c'_{t+n-1}) = R_i.\text{clock}(c'_{t+n-1})$$

となり、式(2), (3)に矛盾する。 \square

以上より、無待機時計合わせプロトコルの同期時間の下界は $\Omega(n)$ である。従って、前節で提案したプロトコルの同期時間はオーダー的に最適である。

6 むすび

本稿では、共有メモリマルチプロセッサシステムにおける無待機時計合わせプロトコルについて考察し

た。同期時間が $O(n)$ の時計合わせプロトコルを提案した。さらに、提案したプロトコルの同期時間がオーダー的に最適であることを示した。今後の課題としては、同期時間が $O(n)$ の自己安定性のある無待機時計合わせプロトコルの考察などが挙げられる。

謝辞

日頃より、分散アルゴリズムについて熱心にご討論いただき本学情報科学センターの片山喜章助手、本学藤原研究室の上田英一郎氏に感謝致します。また、井上智生助手を始めとする本学藤原研究室の諸氏に感謝致します。なお、本研究の一部は、文部省科学研究費補助金(奨励(A)08780279)、財団法人電気通信普及財團の助成による。

参考文献

- [1] L. Lamport and P.M. Melliar-Smith, "Synchronizing Clocks in the Presence of Faults," Journal of the ACM, Vol. 32, No. 1, 1985, pp1-36.
- [2] J.W. Welch and N. Lynch, "A New Fault-Tolerant Algorithm for Clock Synchronization," Information and Computation, Vol. 77, No. 1, 1988, pp1-36.
- [3] D. Dolev, J.Y. Halpern and H. R. Strong, "On the Possibility and Impossibility of Achieving Clock Synchronization," Journal of Computer System Science Vol. 32, No. 2, 1986, pp230-250.
- [4] T. K. Srikanth and S. Toung, "Optimal Clock Synchronization," Journal of the ACM, Vol. 34, No. 3, 1987, pp626-645.
- [5] J. Halpern, B. Simons, R. Strong and D. Dolev, "Fault-Tolerant Clock Synchronization," Proceeding of the 3rd ACM Symposium on Principles of Distributed Computing, 1984, pp89-102.
- [6] M. G. Gouda and T. Herman, "Stabilizing Unison," Information Processing Letters, Vol. 35, 1990, pp171-175.
- [7] A. Arora, S. Dolev and M. Gouda, "Maintaining Digital Clocks in Step," Parallel Processing Letters, Vol. 1 No. 1, 1991, pp11-18.
- [8] S. Dolev and J.L. Welch, "Wait-Free Clock Synchronization," Proceedings of the 12th ACM Symposium on Principles of Distributed Computing, 1993, pp97-108.
- [9] M. Papatriantafilou and P. Tsigas, "On Self-Stabilizing Wait-Free Clock Synchronization," Proceedings of the 4th Scandinavian Workshop on Algorithm Theory(LNCS 824), 1994, pp267-277.

```

1  begin
2    repeat forever do on receipt of a pulse
3      read( $R_j$ );
4      check_a_nap();
5      case work_count of
6        work_count = 0 : wait();
7        work_count = 2n : sort();
8        work_count > 2n : adjust();
9      end
10     count := count + 1;
11     if wait ≠ 0
12       then work_count := work_count + 1;
13     next_read();
14     write( $R_i$ )
15   end.

16  procedure check_a_nap();
17  begin
18    Δcount := count - last[j].my_count;
19    Δcountj :=  $P_j$ .count - last[j].count;
20    dif_Δcount := Δcount - Δcountj;
21    if dif_Δcount > 0 then
22      if  $R_j$ .clock_gen = last[j].clock_gen
23        then invalid[j] :=  $R_j$ .clock_gen;
24      else invalid[j] :=  $R_j$ .clock_gen - 1;
25    if dif_Δcount < 0
26      or  $R_j$ .invalid[i] ≥ clock_gen
27      then local_reset();
28  end;

29  procedure local_reset();
30  begin
31    clock := 0; work_count := 0;
32    adjusted := 0; stay := 0;
33    if wait = 0
34      then clock_gen := clock_gen + 1;
35      wait := n;
36  end;

37  procedure wait();
38  begin
39    wait := wait - 1;
40    position := position + 1(mod n);
41  end;

42  procedure sort();
43  begin
44    for l := 1 to n do
45      if last[i].clock_gen > invalid[i]
46        then now_work_count[l] := last[l].work_count + count
47            - last[l].my_count;
48      else now_work_count[l] := 0;
49    sort  $P_1, \dots, P_n$  (各  $P_l$  の
50      (now_work_count[l], l) の辞書式順)
51      into  $\mathcal{Q} = (P_{q_1}, P_{q_2}, \dots, P_{q_n})$ ;
52    for l := 1 to n do
53      m := 1;
54      while  $P_{q_m} \neq P_l$  do m := m + 1;

56      ranking[l] := m;
57      position := 1;
58      if ranking[i] = 1 then adjusted := 1;
59    end;

60  procedure adjust();
61  begin
62    if adjusted = 1
63      then clock := clock + 1;
64      position := position + 1(mod n);
65    else
66      adjust0();
67      if stay = n
68        then local_reset();
69      if position = ranking[i]
70        then adjusted := 1;
71    end;

72  procedure adjust0();
73  begin
74    if dif_Δcount > 0
75      or  $R_j$ .clock_gen > last[j].clock_gen
76      or clock >  $R_j$ .clock
77    then
78      if clock ≠ 0 then clock := clock + 1;
79      position := position + 1(mod n);
80    else if clock ≠ 0 or clock <  $R_j$ .clock
81      then local_reset();
82    else
83      if  $R_j$ .adjusted = 1
84        then
85          clock :=  $R_j$ .clock + 1;
86          position := position + 1(mod n);
87      else
88        if clock ≠ 0
89          then clock := clock + 1;
90        stay := stay + 1;
91    end;

92  procedure next_read();
93  begin
94    case work_count of
95      work_count = 0 ∨ work_count > 2n :
96        m := 1;
97        while position ≠ ranking[m]
98          do m := m + 1;
99          j := m;
100         1 ≤ work_count ≤ n :
101          j := work_count;
102          n + 1 ≤ work_count ≤ 2n :
103          j := work_count - n;
104    end;
105  end;

```

図 1：同期時間 $O(n)$ の無待機時計
合わせプロトコル