

並列処理における P E 間に渡るゴールの 重みつき参照カウントを用いた管理方式

六沢 一昭 * 市吉 伸行 * 藤 和男 *
吉田かおる * 稲村 雄 * 中島 浩 **

* (財) 新世代コンピュータ技術開発機構 ** 三菱電機(株)

1. はじめに

ゴールの終了の検出と強制終了は、論理型プログラムをマルチプロセッサで並列実行する際の主要課題である。これらは密結合マルチプロセッサでは簡単であるが疎結合マルチプロセッサでは難しい。これは、送信はなされたが先プロセッサにはまだ到着していないようなゴールがネットワーク上に存在しうるためである。

本稿では、疎結合マルチプロセッサであるマルチPSI [1] における Weighted Throw Counting(WTC)方式によるゴールの終了検出と強制終了について述べる。

WTC方式はルートとリーフに重みを持たすことによって管理を行なうものである。この重み付けをGCに適用したものに Weighted Reference Counting方式[3, 4] がある。

2. 計算モデル

マルチPSIは最大64台のPSIをネットワークで結合した並列マシンである。共有メモリではなく各PSI(以下PEと略す)はメッセージ通信によってPE間処理を行なう。同一PE間ではメッセージは送出した順序に到着するが、一定時間内に到着する保証はない。

すべてのゴールは莊園の下で管理される[2]。システムには有限個の莊園があり、莊園IDで識別される。ゴールは新しい莊園を生成することができる。このため莊園はネストすることがあるが、ネストした莊園はゴールと同じように扱うことができるので、本稿では省略する。

負荷分散などのためにゴールが他のPEへ投げ出されることがある。投げ出されてから到着するまで時間がかかるため、ある与えられた時刻においてネットワーク上にゴールが存在する可能性がある(図1)。

ゴールのあるPEには里親がありゴールはその下に属する。里親は1つの莊園に対して1PEに1つだけ存在する。受信したゴールは対応する(同じ莊園IDを持つ)里親に加えられる。対応する里親がない場合は新しく作る。里親は自分の下のゴールの強制終了を行なうことができ、自分の下のすべてのゴールが終了すると消滅する。

3. 問題点

莊園は全ゴールの終了の検出及び強制終了ができなければならない。以下にそれらを行なう際の問題点を述べる。

3-1. 終了の検出

里親は消滅する際その旨を伝えるメッセージ(%terminated)を莊園へ送ることができる。しかしネットワーク上にもゴールがあるかもしれない、たとえ莊園がすべてのPEから%terminatedを受信してもすべてのゴールが終了したとは限らない。あるPEは%terminatedを送信した後、ゴールを再び受信するかもしれない(図2)。

ゴール受信に対して必ず応答を返すならば、ゴールの投げ出しと応答をカウントすることによってゴールの終了を正しく検出することができる[5]。しかしこの方式ではゴ

ール投げ出しと同数の制御(応答)メッセージが必要になってしまう。

3-2. 強制終了

莊園が強制終了を通知するメッセージ(%abortion)をブロードキャストすると、その時点で里親の下にあったゴールの強制終了は可能である。しかしネットワーク上にあったゴールは強制終了できない。

ゴールの強制終了を行なった後も里親が消滅せず、その後受信したゴールを棄ててしまうのならば、ブロードキャストによる強制終了は可能である。

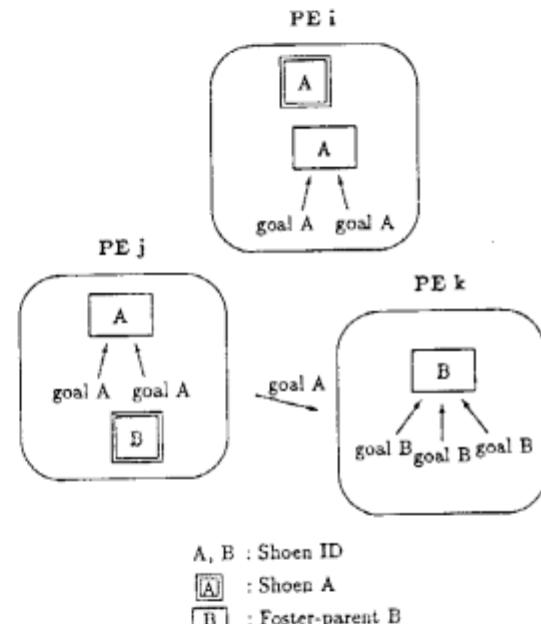


図1. 莊園と里親(Foster-parent), ゴール

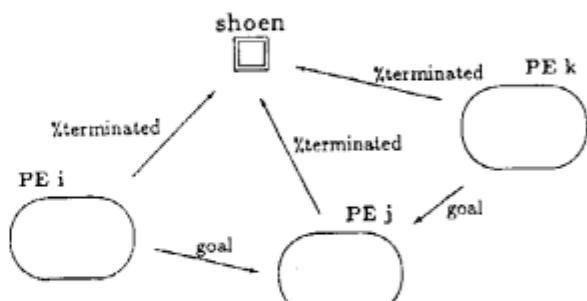


図2. まだゴールがネットワーク上有る。

A Termination Detection and Abortion Scheme for Distributed Processing Systems

Using Weighted Throw Counting

Kazuaki ROKUSAWA, Nobuyuki ICHIYOSHI, Kazuo TAKI, Kaoru YOSHIDA, Yu INAMURA (ICOT),

Hiroshi NAKASHIMA (MELCO)

しかしこの方式には重大な欠点がある。少数の P E にのみ里親が存在した場合多くの %abort が無駄になってしまふ。また莊園 I D の再利用ができない。莊園 I D の再利用を可能にするには、すべてのゴールが強制終了されたことを検出することが必要である。

4. 解決方法

WTC方式による終了の検出と強制終了について述べる。この方式は、制御メッセージの送信は少なく、莊園 I D の再利用も可能である。

4-1. 終了の検出

莊園と里親及びネットワーク上のゴールに“重み”を持たせる。莊園の重みは負の整数、里親及びゴールの重みは正の整数である。そして重みの合計がゼロになるように制御する。この結果、すべてのゴールが終了した時のみ莊園の重みがゼロになる(図3)。

里親はゴールに重みを付けて投げ出す。自分の重みはその分だけ減らす。ゴールを受信するとその重みを対応する里親に加える。対応する里親が存在せず新たに生成した場合は重みの初期値をゴールの重みにセットする。

自分の下にあるゴールがすべて終了すると里親は消滅し莊園へ %terminated を送る。%terminated は消滅した里親が保持していた重みを運ぶ。莊園が %terminated を受信すると、運ばれてきた重みを自分の重みに加える。この操作により重みがゼロになったならば全ゴールの終了が検出される。

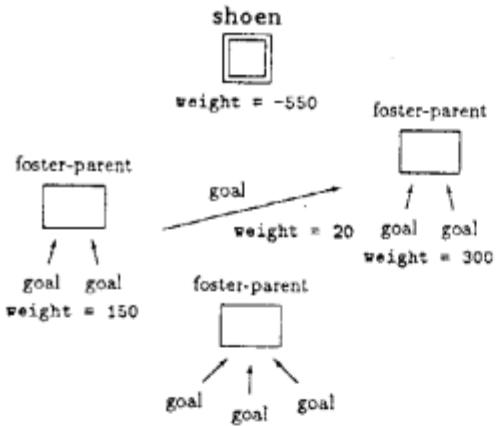


図3. Weighted Throw Counting

4-2. 強制終了

4-1で述べた方法により終了の検出は正しく行われる。従ってすべての里親に %abort が届けばよい。3-2で述べたように %abort をブロードキャストするのでは不十分である。%abort 受信後に生成された里親には %abort が届かないためである。

ここで里親の生成を莊園へ伝えるメッセージ(%ready)を導入する。新たに里親を生成した時 %ready を莊園へ送信する。莊園は %ready を受信すると送信元 P E を記憶する。この記憶は %terminated を受信した際削除する。

強制終了は以下のように %abort を送信することによって行なう。尚 %abort にもゴールと同様重みを付ける。

- ① 記憶している P E へ %abort を送信する。
- ② ①の後 %ready を受信した場合、送信元 P E へ %abort を送信する。

強制終了開始時に莊園が把握していた里親は①によって、②の後に検出された里親は②によって消滅する(図4)。

%abort を受信すると里親は自分の下のゴールをすべて終了させ消滅し %terminated を莊園へ送信する。%terminated は消滅した里親の重みと %abort の重みの合計を運ぶ。里親が既に消滅していた場合は、%abort の重みを莊園へ返却するメッセージ(%return)を莊園へ送る。莊園は %return を受信すると %terminated 受信時と同様に運ばれてきた重みを自分の重みに加える。この操作により重みがゼロになったならば全ゴールの強制終了が検出される。

%abort を1回送信すると1つの里親が消滅しその重みが莊園へ戻る。里親とネットワーク中のゴールが持つ重みの合計は有限なので、有限回の %abort 送信によって強制終了は完了する。

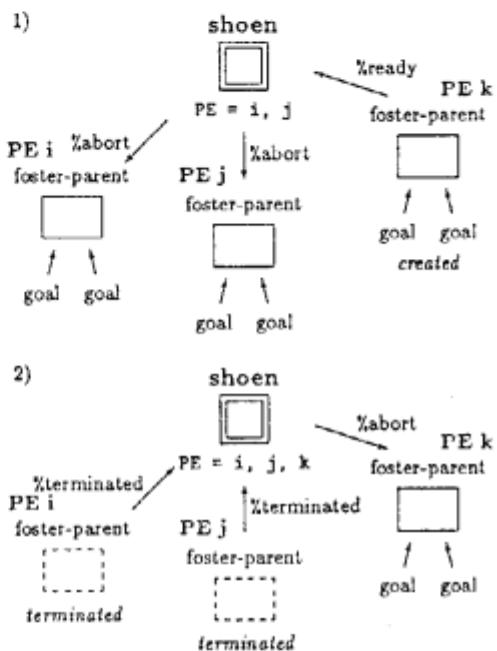


図4. 強制終了

4-3. 重みが1である時の処理

重みが1であるとゴールを投げ出すことができない。この時里親は莊園へ重みの割り付けを要求するメッセージを送信する。このメッセージを受信すると莊園は重みを割り付けるメッセージを返信する。

5. おわりに

WTC方式によるゴールの終了検出と強制終了について述べた。この方式の利点は、「制御メッセージの送信が少なくてすむ。」、「莊園 I D の再利用が可能。」である。

参考文献

- [1] 鹿和男他, "Multi-PSI システムの概要," 第32回情処全大 50-8, 1986.
- [2] 佐藤裕幸他, "並列論理型OS - PIMOS (1)," 第35回情処全大 4D-3, 1987.
- [3] D. I. Bevan, "Distributed garbage collection using reference counting," PARLE, 1987.
- [4] P. Watson and I. Watson, "An efficient garbage collection scheme for parallel computer architectures," PARLE, 1987.
- [5] N. Ichiyoshi et al, "A Distributed Implementation of Flat GHC on the Multi-PSI," ICLP, 1987.