

同期通信用メモリにおけるカウンタとブロッキングの効果

山脇 彰[†](正員) 岩根 雅彦[†](正員)

The Effect of Counter and Blocking for a Communication and Synchronization Memory

Akira YAMAWAKI[†] and Masahiko IWANE[†], Regular Members

[†]九州工業大学工学部, 北九州市

Faculty of Engineering, Kyushu Institute of Technology,
Kitakyushu-shi, 804-8550 Japan

あらまし 同期通信用メモリ TCSM (Tagged Communication and Synchronization Memory) のブロッキング機構について述べ、通信回数を書き込むカウンタの効果と、バスバックオフ機能を適用したブロッキングの効果について明らかにする。実験から、通常のメモリと比較して TCSM はバックオフとカウンタにより実行時間を平均で 39% に抑えた。その速度向上に対して、カウンタによる効果が平均で 58%、バックオフによる効果が平均で 42% となった。

キーワード マルチプロセッサ, 同期通信用メモリ, ブロッキング, カウンタ, バスバックオフ

1. まえがき

マルチプロセッサにおける文や部分文レベルでの並列処理においては、頻繁に発生する同期と通信のオーバーヘッドの削減が重要となる。オーバーヘッドを削減するために、同期と通信を同時に扱う構造化メモリとして、共有メモリにフル・エンティビットを付加した I-structure [1] やそれを拡張してキュー構造をもたせた Q-structure [2], ロックディレクトリを付加したキャッシュ [4], 共有メモリにカウンタを付加した CAM によってエントリの動的な割当てを可能とした TCSM [3] が提案されている。

ここでは、バスバックオフ機能 [5] を適用した TCSM のブロッキング機構について述べる。そして、カウンタとバックオフの効果を明らかにするためにマルチプロセッサ MTA/TCSM II (MultiThread Architecture/TCSM II) [6] で 1 対多通信と相互排除のシミュレーションと実測を行う。

2. 同期通信用メモリのブロッキング

2.1 TCSM の概要 [3]

TCSM は、マルチプロセッサオンチップでのマルチスレッド環境において、相互排除、1 対多通信、条件同期、及びバリア同期を統一的に表現できる同期通信用メモリであり、概念図を図 1 に示す。TCSM は

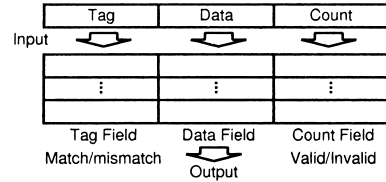


図 1 TCSM の概念図
Fig. 1 Concept of TCSM.

CAM (Content Addressable Memory) で構成され、1 エントリはタグ、データ、カウンタからなる。タグはエントリを識別するためのものであり、カウンタは通信回数を示し、非ゼロ/ゼロによってデータの有効/無効を表す。TCSM は無効なエントリからの読出しと有効なエントリへの書き込みをブロックし生産者-消費者間の同期を保証する。カウンタは有効なエントリの読出し時にデクリメントされる。full/empty ビットとは異なり、TCSM ではカウンタに消費者数を指定できるため、1 対多の同期通信を 1 エントリで実現でき、そのエントリは全消費者による読出し完了後に暗黙的に解放される。

マルチスレッド環境でのタスクはリソースの保護単位であることから、タグをタスク ID と変数名の連結とし、TCSM をタスクごとに保護する。TCSM は同一タスクに属するスレッドまたはスレッドに内在する文や部分文レベルの並列性を抽出したマイクロスレッド間の並列処理に用いられる。

TCSM はマルチプロセッサオンチップ内部で全プロセッサコアから共有され、以下に示す命令で TCSM に対する書き込み、読出し、リセット動作が実行される。STCSM (Store TCSM) 命令は即値のタグと通信回数、レジスタのデータを TCSM に書き込む。LTCSM (Load TCSM) 命令はタグを即値で指定し、レジスタに読み出したデータを格納する。RTCSM (Reset TCSM) 命令はレジスタでマスクパターンを指定し、即値指定のタグで一一致したエントリのカウンタをリセットする。

2.2 バスバックオフ機能の利用

生産者と消費者による書き込みと読出しのブロックをバスバックオフ機能 (BOFF: Bus back-OFF) により実現する。BOFF とは、プロセッサの現在のバス転送を中断させ再実行させる機能であり、BOFF がプロセッサに対してアサートされている間、プロセッサはバスを駆動できない。図 2 に BOFF を適用した TCSM アクセスに関するプロセッサの状態遷移図を

示す．各状態は S_0 (ノンブロック), S_1 (TCSM アクセス完了待ち), S_2 (書込みブロック), S_3 (読出しブロック) からなる．RM (Read Miss) は TCSM に対する読出し失敗, WM (Write Miss) は書込み失敗を表す．RH (Read Hit) は, TCSM に対する読出しが成功しかつそのエントリが解放されたことを, WH (Write Hit) は, TCSM に対する書込みが成功したことを表す．このブロッキング機構を実現するために, 図 3 の BTT (Blocked Tag Table) を組み込んだ．BTT は, ブロック時の状態とそのときのタグを格納するためのテーブルで, タグフィールドは一致検索の CAM である．BTT のエントリがプロセッサに対応しており, BTT 内に格納されるタグは TCSM 内のタグと同一のものである．WB (Write Block flag) は, プロセッサが書込みブロック状態であることを, RB (Read Block flag) は読出しブロック状態であることを表す．これらの初期値は 0 であり, 同時に 1 になることはない．

S_0 のプロセッサが TCSM へアクセスすると S_1 に遷移する．プロセッサが TCSM アクセス中に, WM が発生した場合, BTT のエントリにタグ, WB に 1 がセットされ, RM が発生した場合, BTT のエントリにタグ, RB に 1 がセットされる．このとき BOFF のアサートにより, バス転送が強制的に中断され, TCSM の読出し若しくは書込みに失敗したプロセッサのみがバスを駆動できない状態 (S_2, S_3) になる．BTT に対して, TCSM が RH を発生した場合, タグにヒット

したエントリの WB がリセットされ, TCSM が WH を発生した場合, タグにヒットしたエントリの RB がリセットされる．これにより BOFF がネゲートされ, プロセッサは S_2 若しくは S_3 から S_0 へと遷移する．

3. 実験の準備

3.1 実験環境

対象マシンのマルチプロセッサ MTA/TCSM II [6] は, MTA/TCSM [3] に対して, L1 キャッシュの無効化信号を書き込んだプロセッサ自身には送らないように変更し, TCSM をメモリマップド I/O にした改良機である．また, 高速バリア同期機構への拡張も行っているがここではふれない．MTA/TCSM II は 8 台の 486DX2 [5] からなる単一バス結合共有メモリ型マルチプロセッサであり, 486DX2 内にライトスルーの L1 キャッシュ, マザーボード上にライトバックの L2 キャッシュをもち, TCSM 及びバスアービタは単一バスに接続されている．バスアービタは集中型で, バスの優先度は回転式である．

実験では, プロセッサとスレッドを 1 対 1 に対応させる．シミュレーションでは初期値としてバスの優先順位を ID の最も小さいスレッドが最高で, その後順次回転すると仮定し, 実行時間として, バスサイクル時間の総和とバスサイクルにオーバラップできない命令の実行時間の総和 (T_{no}) を加算し算出した．MTA/TCSM でのシミュレーションによる実行時間 [3] では TCSM アクセスに関して I/O 命令を用いていたが, MTA/TCSM II では mov 命令に変更して算出した．シミュレーションで用いた MTA/TCSM II の基礎データを表 1 に示す．メモリアクセスは最小値である．数値にはバス調停にかかる 2 クロック分を含んでいる．実測では, L1 と L2 ともにオンとし, プログラムのコード及びデータが L1 キャッシュに存在することを確認してから測定した．

3.2 評価プログラム

プログラムの説明で表 2 の記述を用いる．tcsm_w

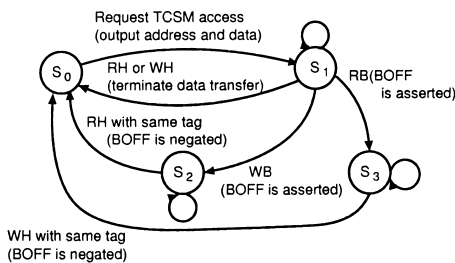


図 2 TCSM アクセスに関するプロセッサの状態遷移図
Fig.2 State transition diagram of processor for TCSM access with bus back off.

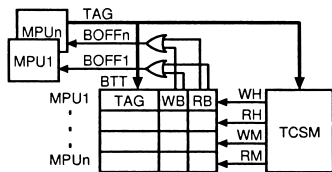


図 3 ブロッキング機構
Fig.3 Blocking mechanism.

表 1 MTA/TCSM II の基礎データ
Table 1 Basic data of MTA/TCSM II.

略称	説明	CLK
W_t	tcsm_w のバスサイクル時間	4
R_t	tcsm_r のバスサイクル時間	4
RB_t	TCSM 読出し失敗のバスサイクル時間	3
W_m	メモリ書込みのバスサイクル時間	4
R_m	メモリ読込みのバスサイクル時間	4(+3)
Bts	Test&Set 命令のバスサイクル時間	9

(+):バースト転送時の追加分

は第1引き数のタグ, 第2引き数のデータ, 第3引き数の通信回数を用いて TCSMへ書き込む. TCSMのブロッキングに BOFF を適用する場合, TCSMの検索と書き込みが同時に行われ, WM時にスレッドはBOFFでブロックされる. BOFFを適用しない場合, tcsm_wはTCSMの検索とTCSMへのデータ書き込みで2回のバスアクセスを発生する. まず, タグを用いてTCSMを検索し WH/WMを調べ, WH時はTCSMにデータを書き込み, WM時は何もしない. 検索と書き込みは不可分動作であり, 検索の結果は戻り値として返されWH時は1を, WM時は0を返す.

tcsm_rは第1引き数のタグを用いてTCSMを読み出し, その結果を第2引き数に返す. BOFFをブロッキングに適用する場合, TCSMの検索と読み出しが同時に行われ, RM時にスレッドはBOFFでブロックされる. BOFFを適用しない場合の動作は, tcsm_wと同様であり, RH時に1を, RM時に0を返す.

評価で使用したプログラムを図4, 図5に示す. 両図において, TCSM_B, TCSM_NBはTCSMを用いて実装したプログラムであり, 前者はBOFFによるブロックを行い, 後者はBOFFによるブロックを行わない. MEMORYはメモリを用いて実装したプログラム

である. 図4のプログラムにおいて, TCSM_Bは生成データに対する消費者の読み出しブロック, MEMORYはデータ生成フラグ, TCSM_NBは tcsm_r の戻り値チェックによって通信データの正しさを保証する. 図5のプログラムにおけるロックに関して, TCSM_BはTCSMの読み出しブロック, TCSM_NBは tcsm_r を用いたピジーウェイト, MEMORYはTEST&SET命令によるスピンロックで実現している.

4. シミュレーション [3]

4.1 1対多通信

図4のプログラムをスレッド ($p = 2 \sim 8$) とし, 生産者のスレッドIDが最も小さい状態でそれぞれをプロセッサに割り当てた. TCSM_Bの実行時間 (T_b), MEMORYの実行時間 (T_m) は次式となる,

$$T_b = 4 \cdot p \tag{1}$$

$$T_m = 21 \cdot p - 13 \tag{2}$$

TCSM_NBは, TCSM_Bに対し消費者の tcsm_r による戻り値チェック分 ($p - 1$ 回), 実行時間が増加する. したがって, 実行時間 (T_{nb}) は次式となる.

$$\begin{aligned} T_{nb} &= W_t + (p - 1)(R_t + R_r) + T_{no} \\ &= 8 \cdot p - 4 \end{aligned} \tag{3}$$

シミュレーション結果を表3に示す.

MEMORYはフラグとデータが別であるため, データとフラグとしてのカウンタを同時に読み込めるTCSM_NBと比較し, $p - 1$ 回バスアクセスが多い. TCSM_Bは, ブロッキングにより tcsm_r の戻り値を調べる必要がないため, TCSM_NBと比べ $p - 1$ 回のTCSMアクセスを削減できる. MEMORYと比較して, TCSM_Bは実行時間を22%に, TCSM_NBは56%に削減している. TCSM_BのMEMORYに対する速度向上の内訳はカウンタで78%, バックオフで22%である.

4.2 相互排除のシミュレーション

図5のプログラムをスレッド ($p = 2 \sim 8$) とし, それ

表2 TCSM アクセスに関する操作の記述
Table 2 Descriptions for TCSM access.

表記	説明
WH=tcsm_w(tag,data,count)	TCSM 書き込み動作を実行
RH=tcsm_r(tag,&data)	TCSM 読み出し動作を実行

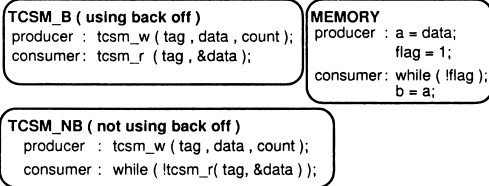


図4 1対多通信のプログラム
Fig. 4 Programs of multicast.

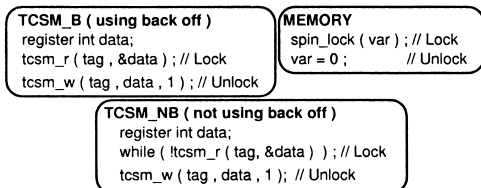


図5 相互排除のプログラム
Fig. 5 Programs of mutual exclusion.

表3 シミュレーション結果 (1対多通信)
Table 3 Result of the simulation. (multi cast)

	$p = 2$	3	4	5	6	7	8
T_b	8	12	16	20	24	28	32
T_{nb}	12	20	28	36	44	52	60
T_m	29	50	71	92	113	134	155

表4 シミュレーション結果(相互排除)

Table 4 Result of the simulation. (mutual exclusion)

	$p=2$	3	4	5	6	7	8
T_b	19	33	50	70	93	119	148
T_{nb}	24	48	80	120	168	224	288
T_m	35	66	106	155	213	280	356

それをプロセッサに割り当てた．TCSM.BとMEMORYの実行時間は次式となる．

$$T_{tb} = \frac{3}{2} \cdot p \cdot (p-1) + 8 \cdot p \quad (4)$$

$$T_m = \frac{1}{2} \cdot p \cdot (9p + 17) \quad (5)$$

TCSM.NBはTCSM.Bに対し、tscm.rによる戻り値チェック分($p-1$ 回)だけバスアクセス回数が増加する．したがって、実行時間は次式となる．

$$\begin{aligned} T_{nb} &= p \cdot R_t + p \cdot (p-1)R_t + p \cdot W_t + T_{no} \\ &= 4 \cdot p^2 + 4 \cdot p \quad (6) \end{aligned}$$

シミュレーション結果を表4に示す．

TCSM.BはTCSM.NBでのtscm.rの戻り値チェック分だけ実行時間を抑えることができた．MEMORYが最も悪い結果となったのは式(4),(5),(6)よりメモリのlock付きTEST&SET命令がバスを9クロック使用するためである．MEMORYと比較してTCSM.Bは実行時間を56%に、TCSM.NBは22%に抑えた．TCSM.BのMEMORYに対する速度向上の内訳は、カウンタで38%、バックオフで62%である．

5. MTA/TCSMでの実測

前章での結果を得て、BOFFによるブロッキングを適用したTCSMをMTA/TCSMで評価した．各プログラムの実行に関する条件はシミュレーションと同じである．1対多通信の結果を表5に、相互排除の結果を表6に示す．表5の p は消費者数であり、数値は生産者がデータを生成してから最後の消費者がデータを読み出すまでにかかったクロック数である．表6の p は全スレッド数を表しており、2~8台までスレッドを変化させている．数値は全スレッドがプログラムの実行を開始してから完了するまでにかかったクロック数である．相互排除において、MEMORYの値がシミュレーションの結果と異なっている．lock付きTEST&SET命令でのメモリアクセスはL2キャッシュにキャッシングされず主メモリ(DRAM)アクセスとなる[5]．したがって、DRAMのリフレッシュによるレイテンシの影響と考える．双方ともTCSM.B

表5 1対多通信の実測結果

Table 5 Actual survey result of multi cast.

	$p=1$	2	3	4	5	6	7
TCSM.B	8	12	16	20	24	28	32
MEMORY	29	50	71	92	113	134	155

表6 相互排除の実測結果

Table 6 Actual survey result of mutual exclusion.

	$p=2$	3	4	5	6	7	8
TCSM.B	19	33	50	70	93	119	148
MEMORY	35	66	115	166	228	310	389

の方が良い結果となっており、1対多通信及び相互排除におけるカウンタとブロッキングの効果を実機で確認できた．

6. むすび

同期通信用メモリTCSMのブロッキング機構について述べ、その機構を搭載したマルチプロセッサMTA/TCSMで1対多通信と相互排除の実験を行った．通常のメモリと比較して、TCSMはバックオフとカウンタにより実行時間を22~56%に抑えた．その速度向上に対して、カウンタによる効果が38~78%、バックオフによる効果が22~62%であり、それぞれの効果についても確認できた．今後は、MTA/TCSM IIでの高速バリア同期機構に対する評価と一般的なプログラムによる評価を行う．

謝辞 日ごろからお世話になっている春日工作所(株)に深謝します．

文 献

- [1] A. Nikhil and R.S. Nikhil, "I-structure: Data structures for parallel computing," Trans. Prog. Lang. and Sys. ACM, vol.11, no.4, pp.598-639, Oct. 1989.
- [2] 佐藤三久, 児玉祐悦, 坂井修一, 山口喜教, "並列計算機EM-4における分散データ構造を用いたマルチスレッドプログラミング;" 情処学 ARC 研報, ARC-92, 1992.
- [3] 岩根雅彦, 山脇 彰, 田中 誠, "マルチプロセッサオンチップにおけるCAMを用いた同期通信用メモリ;" 信学論(D-I), vol.J83-D-I, no.3, pp.317-328, March 2000.
- [4] T. Tarui, T. Nakagawa, N. Ido, M. Asaie, and M. Sugie, "Evaluation of the lock mechanism in a snooping cache," International Conference on Supercomputing 92, pp.53-62, July 1992.
- [5] Intel Corp, "インテル 486 マイクロプロセッサデータブック;" Intel Corp, 1992.
- [6] 山脇 彰, 岩根雅彦, "バスバックオフ機能を用いた同期通信用メモリの制御方式;" 情処学第61回全国大会予稿集, 6D-7, pp.77-78, 2000.
(平成12年12月25日受付, 13年3月29日再受付)