

DSM を用いたプロセス間通信における優先度制御 A Priority Control for Inter-Process Communication with DSM

田中 晶†
Akira Tanaka

1. はじめに

我々は、物理的に離れたノード間でネットワークワイドメモリ空間を共有する DSM (Distributed Shared Memory) 機構を用いて、ノード A の AP (Application (図 1,)) 内のプロセス (以下 AP で総称する) が自ノードの DSM にメッセージを書き込むと (図 1, (a)), ハードウェアによりノード B の同一番地の DSM にメッセージが自動的にコピーされる (図 1, (b)(f)(a')), 高速通信方式を開発してきた[1][2]。この方式は、DSM 書込みとノード間メッセージ転送をオーバーラップして実行するので、小さなレイテンシで転送が行える。ノード間ネットワークではソフトウェア負荷の少ないハードウェア組込みの通信手段も注目され、DSM 通信機構もその小さなレイテンシから有効である。

しかしながら、DSM 書込み処理速度よりネットワーク転送速度が小さいと、ネットワークが隘路になり大規模メッセージが転送し終わるまで次の送信ができず、ノード間制御情報のような実時間性を要する小規模メッセージの遅延を引き起こす可能性がある。そこで、大規模メッセージのフラグメンテーションにより、この課題を解決する優先度制御を提案する。

2. 基本概念と制御方法

小規模メッセージの実時間通信を確保するために、ルータによるフラグメンテーション技術を参考に、次のようなスケジューリングとバッファリングを組み合わせた方式を検討した。本提案で主に必要な構成要素は図 1 に “提案” として示すローカルバッファとその管理機能である。

<スケジューリング> 小規模メッセージの優先度を大規模メッセージの優先度より上げ、大規模メッセージの転送中に小規模メッセージ転送要求が発生するとその処理を中断し、小規模メッセージの転送を行った後、再び大規模メッセージの転送を開始する。

<バッファリング> OS は大規模メッセージに対して、メッセージ毎に所定の規模単位 (一括転送単位、図 1 の *) に区切ってローカルバッファ (BA: Buffer Area) を割り当て、AP (図 1,) は BA に大規模メッセージを書き込む (図 1, (c))。OS は 1 つの一括転送単位が一杯になる毎に一括転送単位を DSM へコピーし、並行して DSM コピーハードウェア (DMC: Distributed Memory Coupler) が相手ノードへ転送する。また、OS は次の転送のため一括転送単位を更新する (図 1, 注)。小規模メッセージは、従来通り AP から直接 DSM に書き込まれ、DMC によりノード間転送がすぐに開始される。

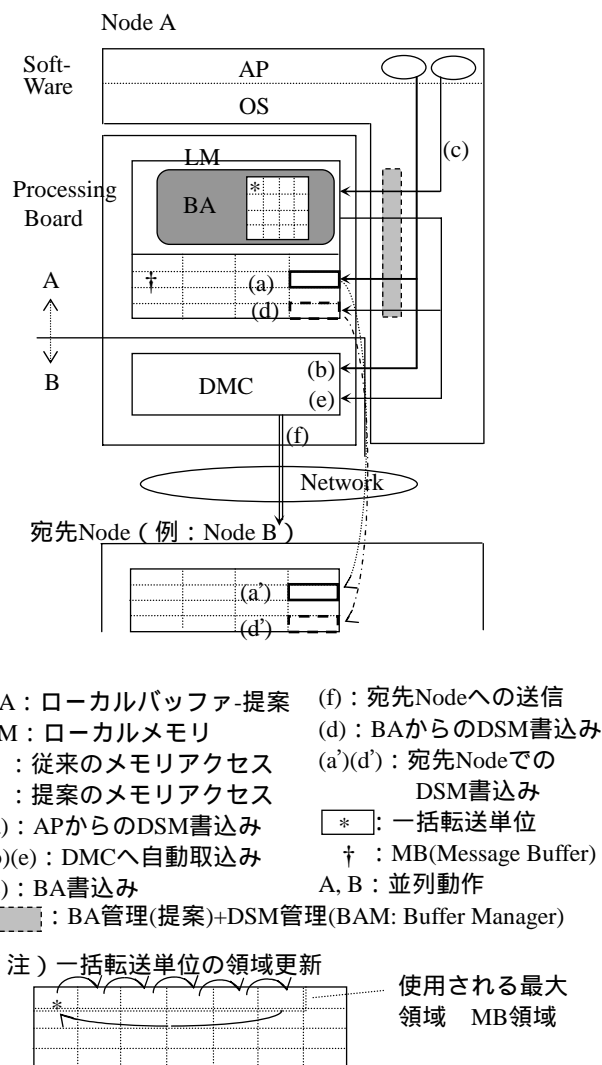


図1 DSM通信の動作と提案機能

このため、DSM の領域 (MB: Message Buffer) 割当てや MB 空塞判定を行う OS の DSM 管理機能を、その対象を BA まで拡張し (以降 BAM (Buffer Manager) と呼称)、バッファリングの際の一括転送単位の割当て、一括転送単位空塞判定、及び BA から DSM へのコピーを行う。

3. 提案手法による DSM 通信の動作

AP から発生したメッセージが相手ノードへ送信される手順を図 2-2 で示す ([] は図 2 の記号)。

- メッセージ送信開始時、AP は DSM 書込みのために MB 確保要求と共にメッセージ規模を OS (BAM) に送る。

†日本電信電話 (株) NTT 未来ねっと研究所,
NTT Network Innovation Laboratories, NTT Corporation

- ii. BAMはそのメッセージ規模を元に一括転送単位の規模を定める。優先度の低い大規模メッセージを送信するAP〔Px〕では一括転送単位の規模>0となり、BAMはDSMの番地の代わりにBAの番地をAPに返す。優先度の高い小規模の実時間メッセージを送信するAP〔Pa〕では一括転送単位の規模は0となり、BAMはDSM内に確保するMBの番地をAPに返す(以下、一括転送単位の規模>0の送信要求が先に生じた場合を述べる)。
- iii. APがメッセージをiiで番地指定されたBA内の一括転送単位に書き込む〔A〕。BAへの転送量が指定された一括転送単位の規模に達するとDMCがそれを検出し〔 〕BAMに通知する。
- iv. 小規模メッセージの送信要求が無ければ、BAMはBAからDSMへコピーを行った後、一括転送単位毎に番地を更新しながら、APに残りのメッセージをBAに書き込むよう指示し、iiiとivを繰り返す。
- v. iiiにおいて、小規模メッセージの送信、即ち一括転送単位=0の優先度の高い送信要求があれば〔Pa〕、BAMは一括転送単位一杯の際に大規模メッセージのBAからDSMへのコピーを保留し、小規模メッセージにDSMのMBを割当て、書き込みを指示する〔 〕。小規模メッセージ送信APはメッセージをDSMに書き込み、DMCが並行してノード間転送する〔B〕。DSM書き込み終了後〔 〕、BAに保留していた大規模メッセージを、BAMがBAから該当DSM領域へ転送する〔C〕。DSM領域への転送終了後〔 〕、Pxのメッセージの残り、或いは、Pxの送信が完了した場合は他のメッセージの送信を引き続き行う〔D〕。
- vi. APが送信を終了するとAPがBAMに通知し、BAMがBA内の一括転送単位を解放する。

即ち、小規模メッセージ送信要求の際、大規模メッセージの送信中でなければ、PaへのMB割当て、DSM書き込み〔B〕とノード間転送のみが直ちに行われ、大規模メッセージ送信中であれば一括転送単位一杯を契機にその処理を中断して、小規模メッセージを送信することになる。尚、ivでの書き込み番地の更新は、DSM転送が完了する前に一括転送単位が上書きされるのを防ぐためである。従って、1ノード向けにバッファリング用に確保されるメモリ規模は、最大で、本来のDSM通信機構が1回で送れる規模、即ちせいぜい1ノード向けのMB[1]領域の規模程度を、図1の注に示すように繰り返し利用する。

4. 本優先度制御の効果とまとめ

図2-2に示すように、大規模メッセージPxをバッファリングし、さらに、一括転送単位毎にスケジューリングにより分割することで、実時間性の高いPaが割込む形で送信される。その結果としてネットワーク側の待ちが回避され、従来図2-1の*1であった実時間メッセージPaの送信を、図2-2の*2のように大規模メッセージに先行して行える。そして次の利点を持つ。(a)一括転送単位の規模はAPが通知するメッセージ規模に基づきBAMが定めるため、間接的ながらAP自身がその実時間性により優先度制御が可能である。これらは、タスク管理に類する複雑な制御無しに比較的簡易に実現できる。(b)APは単にBAMの返す番地にメッセージを書き込むだけであり、実

時間性や規模によらずAPインターフェースは共通で良い。(c)送信の優先度、即ち、転送保留の有無及び送信順序を、唯一のパラメータである一括転送単位の規模のみにより設定できる。

本提案は、送信処理速度と比ネットワーク速度が遅い場合や実時間性の低いメッセージ規模が大きい場合に効果が高く、一方、ネットワークでの遅延が殆ど無い場合は効果が得にくい。このような特徴を活かせるAP、例えば、システムバックアップ情報の送信と混在して制御情報を即時に送る場合等に有効である。

DSMによる高速通信方式において、既存の主要構成要素であるDMCを改造することなく、OSの機能を少しばかり拡張することで、大規模メッセージが混在しても、制御情報等の小規模メッセージの実時間性を確保できる技術を明らかにした。

【参考文献】

- [1] 田中 晶, 山田茂樹, 田中 聡: ネットワーク分散処理ノードアーキテクチャ MESCARのメモリ間複製機構の設計と評価, 電子情報通信学会論文誌 B, Vol.J86-B, No.2, pp.148-161 (2003).
- [2] 山田茂樹, 田中 聡, 田中 晶, 向井良: メモリ結合型分散処理ノードアーキテクチャ MESCARのインプリメンテーションと性能評価, 電子情報通信学会論文誌 D-I, Vol.J83-D-I, No.4, pp.418-429 (2000).

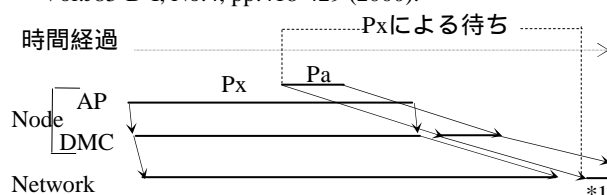
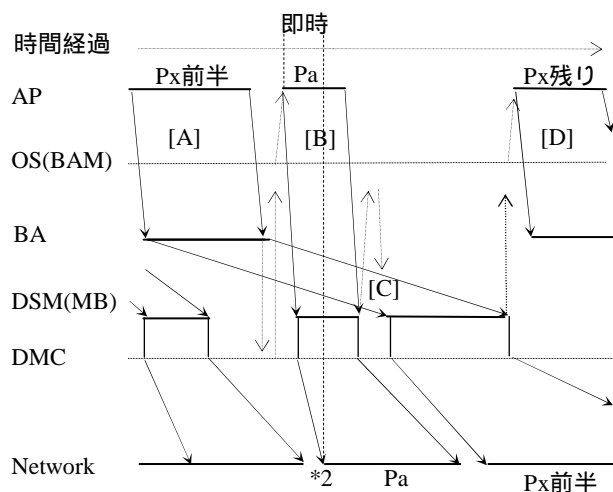


図2-1 従来手法



- ：BA書き込みが一括転送単位に達したことを検出。
- ：高優先度メッセージがあれば書き込み開始。
- ：PaのDSM書き込み終了を検出。
- ：Px前半のBA->MB転送を開始。
- ：Px前半のBA->MB転送終了を検出。
- ：Px残り或いは次メッセージ転送を開始。

図2-2 スケジューリングとバッファリング

Pa: 実時間メッセージ Px: 大規模メッセージ

図2 提案手法導入の効果