

## 汎用超並列オペレーティングシステム SSS-CORE のメモリベース通信機能

5B-3

松本 尚 平木 敬

東京大学 大学院理学系研究科 情報科学専攻\*

## 1 はじめに

汎用超並列オペレーティングシステム SSS-CORE[1] は並列アプリケーションと協調動作することで、効率を極力落すことなくマルチユーザ/マルチジョブの汎用環境を実現する分散メモリ型並列計算機およびワークステーションクラス環境(NOW: Network of Workstations)を対象とした汎用オペレーティングシステム(汎用OS)である。SSS-COREはシステムの資源管理に階層性を導入して資源管理の効率化を行うことにより、スケラビリティつまり超並列超分散計算環境に対応している。ユーザの並列アプリケーションの効率の良い実行のためには、もちろん従来 SSS-CORE で主張していたユーザ/カーネルの協調資源割当や資源管理効率化によるカーネルコストの削減も重要である。しかし、第一義的にはユーザモードにおいてノード間における通信と同期をいかに高速に実現するかが最大の鍵である。本稿では特殊な通信同期ハードウェアを仮定しないNOW環境においても、高速なユーザ通信/ユーザ同期を提供するメモリベース通信機能の基本方針と実装方式の概略を示す。

## 2 高速ユーザ通信同期実現の問題点

以下に高速ユーザ通信同期実現のための障害となってきた問題点を列挙する。

## 2.1 ユーザ/カーネルおよびコンテキスト切替コスト

汎用環境を実現しようとしているため高速ユーザ通信同期であってもOSによるプロテクションを廃止できない。現在の多くの分散メモリ型高並列計算機では各ノードでシングルノード用のOSを動かして、従来OSの保護された通信同期機能を利用している。このため、非常に大きなユーザ/カーネル切替およびコンテキスト切替のオーバーヘッドが通信と同期のたびに必要となる。

## 2.2 メッセージキュー操作のオーバーヘッド

多くの分散メモリ型並列計算機やNOWでは他ノードから飛んで来たメッセージパケットを少数個(通常1個)の通信用デバイスで受信して、とりあえずメモリ内の受信バッファ領域に格納する。ここまでの処理はハードウェアの仕様で固定されており通常回避不可能である。その後、多くのシステムでは以下のような手順となる。OSがそのメッセージをプロセス毎に仕分けしてハードウェアの専用バッファ領域から移動し、メッセージを待っていたユーザプロセスを起動状態にする。起動されたユーザプロセスはメッセージの内容を読んで処理内容を選択するために移動して、さらに最終的な処理を行うための領域に移動して、依頼された処理を行う。このようにパケットデータの移動や解読が場所を返して繰り返されて効率が悪い。

## 2.3 ノード間コネクション保持のオーバーヘッド

前記のメッセージキュー操作のオーバーヘッドを緩和するために、他ノードのユーザメモリ空間を比較的低コストで操作する遠隔メモリアクセスを分散メモリ型並列計算機やNOWでも取り入れつつある。しかし、この場合でもノード間のメモリのマッピングをユーザもしくはOSが管理保持しなくてはならないので、超並列超分散環境で通信同期するノード数および頻度が大きくなるとコストが増大する。

## 2.4 マルチキャスト通信のオーバーヘッド

コネクション保持のコストとも関連するが、対象ノード数が増大すると1-to-1ベースの通信手段の繰り返しで実現される1-to-many通信のオーバーヘッドが非常に大きい。特殊なブロードキャスト通信手段がシステムに存在しても、Acknowledge(Ack)が必要な場面ではAckの収集に大きなコストがかかってしまう。

## 2.5 キャッシュ/TLBのポリューション

メッセージ駆動型の実行形態やデータ駆動型の実行形態や通常の高性能マイクロプロセッサを利用した並列計算機におけるActive Message[4]流の実行形態では、処理が本来持つ局所性の利用が難しい。汎用環境では無関係なジョブの実行が細粒度で混じり合うので、キャッシュやTLBのポリューションの度合がより一層悪化し、性能が低下してしまう。現在の局所性の利用が大前提となっている高性能マイクロプロセッサを要素プロセッサとして使用する限り、これは大きな欠点である。

## 3 メモリベース通信機能の特徴

前記の問題点が認識されれば解決策は簡単であり、それは共有メモリプログラミングモデルに基づく高機能遠隔メモリアクセス、つまりメモリベース通信(同期)機能である。SSS-COREにおけるメモリベース通信機能とはMemory-Based Processor(MBP)[2]と共に提案された高機能分散共有メモリシステム(SMS: Strategic Memory System)上の各種メモリベース通信機能やメモリベース同期機能[3]を、NOWや分散メモリ型並列計算機に可能な限り低オーバーヘッドのソフトウェアにより実装したものに他ならない。以下、どのようにして前記の問題点が解決されるかを簡単に説明することにより、メモリベース通信(同期)機能の特徴を述べる。

- 単にユーザ/カーネルのモード切替トラップや外部割込ハンドラへの反応時間だけで見れば、高性能マイクロプロセッサのオーバーヘッドは数クロック~数十クロック程度である。そこで付随するオーバーヘッドを最小限に抑えるために以下のような実装を行う。ユーザレベルで使用するノード間通信同期用システムコールや割込ハンドラは、OSの他のI/O関係のシステムコールや割込ハンドラと分離して実装し、通信同期にとって余分なチェック等を全廃し、カーネル権限で実行されるコードおよびユーザ/カーネル切替に際して退避復旧するプロセッサコンテキストを最小限にして、実装される。この通信同期用システムコールや割込ハンドラではアドレス空間の切替を行わない。
- OSのプロテクションの下でバッファのコピー回数やソフトウェアオーバーヘッドを抑える最良の方法は、メモリ管理機構と協調動作する遠隔メモリアクセス用のハードウェア(つまりMBPのような物)を使用することである。この遠隔メモリアクセス用ハードウェアを仮定しない場合は、通信パケット内に通信相手先のアドレス空間内の目的アドレスを格納しておき、1.の実装方針に従って作られたメッセージ受信用の割込ハンドラが直接相手先のユーザ空間内のメモリ領域に通信パケット内のデータを受信バッファ領域から格納する。プロテクションはパケットの発信元と受信先の少なくともどちらか一方でケイバリティチェックを行うことで実現できる。SSS-COREの現在の実装ではパスワード(32bit数値)を利用したケイバリティチェックを受信先の割込ハンドラで行っている。
- 動的な通信はコネクションも動的に接続・解除されるので、コネクション保持の問題はあまり生じない。しかし、静的(コンパイル時)に、通信相手や送信先アドレスやフェッチデータのアドレスが判る場合には、通信相手や転送先アドレスといったコネクションに関する静的情報を利用した方が、動的なコネクションを張るオーバーヘッドが回避できるので効率が良い。並列実行に非常に適したデータパラレル(SPM)の数値計算アプリケーションでは、論理アドレス空間を共通(ただしメモリの実体は別)にすることが非常に有効である。論理アドレス空間を共通にすることにより、遠隔メモリアクセスのためのコネクション情報はプログラムコード内に操作対象の論理アドレス(と静的に判明する論理プロセッサ)として格

\*Memory-Based Communication Facilities of the General-Purpose Massively-Parallel Operating System: SSS-CORE. Takashi MATSUMOTO and Kei HIRAKI, University of Tokyo, tm@is.s.u-tokyo.ac.jp

納可能になり、実行時には論理プロセッサと物理プロセッサ（物理ノード）の対応を取るだけで通信同期が可能である。通信同期要求パケットの受信割込ハンドラが対象論理プロセッサの論理アドレスを解決<sup>1</sup>してメモリ操作を行う。

- MBPによるSMSの大きな特徴に階層マルチキャストとAckコンバイニングによるupdateベースのキャッシュの実現がある。このupdateベースのキャッシュシステムは上記のSPMDプログラムの例と同様に論理アドレス空間を各ノードで共通に取り、論理アドレスが同じキャッシュ用のページを各ノードに用意して、階層マルチキャストやAckコンバイニングをパケット受信割込ハンドラでエミュレートすることで実現できる。基本的に大規模なマルチキャストは階層構造を利用して行われ、Ackのコンバイニングも階層構造を利用しない限り効率化できないので、このupdateキャッシュの実現方式はNOWであれ、分散メモリ並列計算機であれ、変わりはない。そこで、NOWにおいてはネットワーク内に階層的な経路を静的に見出し、それを利用して実現される。
- ナイーブなメッセージ駆動や要求駆動などのパケット到着駆動による実行形態を基本実行形態として選択するには現在の高性能マルチプロセッサは処理の空間的・時間的局所性を期待しすぎている。また、局所性の利用は、コンピュータの高速化にとって最も重要な事項であるから、不可避である。しかし、処理したいデータセットの大部分が他のノードにあるような場合は逆に通信によってコントロールを移動した方が局所性が抽出できる。この場合においても、現在実行中のジョブからCPUを奪って実行したのでは、現在実行中のジョブの局所性の利用を妨げ、スケジューリングの公平性からも望ましくない。このためSSS-COREでは、SMSのMemory-Based Signalと同様に、ユーザレベルのランキューにパケットとして運ばれてきた実行スレッドを直接格納することで、キャッシュポリューションや不公平性を回避する。このユーザレベルランキューの所有者であるジョブが実プロセッサにスケジューリングされている時しかこのキューからスレッドが起動されることはない。このため関係のないジョブのCPU時間の消費はキューへの登録操作（データ量が多い場合は実際のデータ移動はDMA転送）のみであり最小限である。さらに、このスレッド起動はアドレス空間がすでに切り替わっており、ユーザレベルで実現されるので、極めて低コストである。

また、遠隔メモリリードや各種メモリベース同期機能は要求元からの遠隔メモリアクセス要求パケットの送信、受信時に割込ハンドラ内で目的クラスタ（ノード）内の目的タスクの対象アドレスの操作（同期処理の種類によっては不可分操作）および要求元への返り値（返り値）の返送、要求元での受信割込ハンドラによる返り値の返り値アドレスへの書き込み（Ackの場合はカウントアップやカウントダウン）で実現される。

ノード間通信同期に関わる低レベルの各種同期情報はすべてユーザレベルのフラグに反映され、このフラグを利用してSnoopy Spin wait[5]によって同期が行われる。遠隔メモリアクセスの順序モデルには緩和されたメモリアクセスモデルを利用することで性能を向上させる。遠隔メモリアクセスの順序管理はAckベースのメモリバリアで行われるが、可能な限りAckを省略して済ませるように実装を行う。また、実装上の問題としては、EtherやFastEtherのようにデバイスレベルでパケットが消失する可能性がある通信ハードウェアに対応するために、パケットに送信ノード毎にシリアルナンバーをつけて管理している。紙面に限りがあるため、メモリベース通信機能およびその実装方式の詳細については別稿で報告する予定である。

#### 4 メモリベース通信機能に適したプログラミングモデル

前出のupdate系書き込みと遠隔メモリ書き込みアクセスは、プログラミングモデル上はすべてのノードにおいて単なるメモリアクセスとして記述されるが、NOW上のSSS-COREにおける実行コードでは、これらのアクセスは遠隔メモリアクセスのための一連のコードシーケンスとしてコンパイルされることで効率良く実現される。このためupdate系書き込み（updateページへの書き込み）や遠隔メモリ書き込み（home\_onlyページへの書き込み）が静的に解析可能なプログラミ

ングモデルを採用する必要がある。さらに、update系書き込みの転送先リストも静的に決定されている方が実行時の効率が良い（SMSでもupdateのコピーページは回収されない）。また、頻度が少なく一部しか外部から参照されないページはコピーページを作らずに遠隔メモリ（home\_onlyページからの）読み出し（またはプリフェッチ）で対応した方が、性能上もメモリ資源節約上も効率が良い。

大きな粒度でノード間でデータ交換すれば良いデータはinvalidateページに割り当て、invalidateページはIVY[6]と同様のページ管理機構を利用した無効化プロトコルベースの仮想共有メモリとして実現される。なお、update、invalidate、home\_onlyの属性は共有メモリ空間のある一時点においては一意でなくてはならない。つまり、同一の共有メモリ領域に対して、あるタスク<sup>2</sup>ではコピーを持つupdate領域で、他のあるタスクからはコピーを持たないhome\_only領域というようなメモリ領域を同一論理アドレスで定義することはできない。ただし、OSによって同一論理アドレス領域の属性を一斉に変更することは禁止しない。この他にlocalページ（node.localではなくprocessor.localの意）というメモリ領域に対する論理的な属性値が存在し、この属性値の領域はたとえプログラム上の論理アドレスが同一であり同一クラスタ内であってもプロセッサごとに別の独立した実アドレス領域に割り当てられる（つまり別タスクとなる）。このlocal属性は集共有メモリマルチプロセッサのノード（クラスタ）を持つシステムにおいて、SPMD型のプログラムを効率良く実行（ローカル変数領域の同一論理アドレス上の確保）するために新たにSSS-COREのメモリシステムに付加された属性である。この属性のメモリ領域を持つタスクは同一時点では多くとも一つの実行中のスレッド（実プロセッサ）しか持たない。

#### 5 おわりに

現在、ワークステーションクラスタ版SSS-COREの開発はSun Microsystems社のSPARCstation 10またはSPARCstation 20をEthernetで接続した環境で動作している。WS単体上の機能はすでに充実しておりメモリ保護、タイムシェアリング、プロセス管理、UDP通信、TCP/IP通信、キー入力、画面出力、外部UNIXワークステーションからのプログラムロード/実行等が行える。本稿で述べたメモリベース通信機能も一番ベースとなる遠隔メモリ書き込みに関してはすでに実装されており、これを利用した並列レイトレーシングや並列メンデルブロー集合表示のデモプログラムがSSS-CORE上で動作する。これからメモリベース通信機能を拡充していくと共に、ワークステーション間に跨る資源管理方式やユーザレベルのカーネル協調ランタイム/ライブラリも早急に開発実装していく予定である。また、実装と並行して各種性能データを採取する予定である。

#### 謝辞

本研究は情報処理振興事業会（IPA）が実施している独創的情報技術育成事業の一環として行なった。ワークステーション版SSS-COREの共同開発者である有限会社アックスの駒嵐丈人氏と竹岡尚三氏に感謝いたします。また、本研究に関して有益な議論をいただいた有限会社アックスの湯原茂氏に深謝します。

#### 参考文献

- 松本 尚, 平木 敬: 汎用並列オペレーティングシステム SSS-CORE の資源管理方式日本ソフトウェア科学会第11回大会論文集, pp.13-16 (October 1994).
- 松本 尚, 平木 敬: Memory-Based Processor による分散共有メモリ、並列処理シンポジウム JSPP '93 論文集, pp.245-252 (May 1993).
- 松本 尚, 平木 敬: キャッシュインジェクションとメモリベース同期機構の高速化。計算機アーキテクチャ研究会報告 No.101-15, 情報処理学会, pp.113-120 (August 1993).
- T. von Eicken, D. E. Culler et al.: Active Messages: A Mechanism for Integrated Communication and Computation. *Proc. 19th Int. Symp. on Computer Architecture*, pp.256-266 (May 1992).
- 松本 尚: マルチプロセッサ上の同期機構とプロセッサスケジューリングに関する考察。計算機アーキテクチャ研究会報告 No.79-1, 情報処理学会, pp.1-8 (November 1989).
- K. Li: IVY: A Shared Virtual Memory System for Parallel Computing. *Proc. 1988 Int. Conf. on Parallel Processing*, St. Charls, IL, pp.94-101 (August 1988).

<sup>1</sup>解決と言っても通常TLBを多くとも1本セットするだけである。

<sup>2</sup>メモリ空間の割り当て単位であり、同一タスクに属するプロセッサは同一ノード内に存在しページテーブルを完全に共有している。