

流れ場シミュレーションにおける PVM の使用法

PVM Implementation of Parallel Computation in NST

松田進也*・大島まり*・谷口伸行*

Shinya MATSUDA, Marie OSHIMA and Nobuyuki TANIGUCHI

1. PVM について

PVM¹⁾ (Parallel Virtual Machine) とは、ネットワークによって接続されたコンピュータ間で、メッセージパッキングを用いて並列・分散処理環境を構築するツールです。超並列マシンのようにコンパイラが勝手に並列化してくれたり、プログラム中に並列化指示を書いておけばすむようなものとは違い、設計の段階から並列化を意識して行わなければならない、後述するようにデバッグも難しくなりますが、ネットワークでつながってさえいれば UNIX から Windows までどんな機種でも使えること、並列化の自由度が高いことなどがメリットとして挙げられます。この PVM を用いて行ったトルクコンバータ内流れの計算を例に、PVM の具体的な使用法を説明しようと思います。

並列化には大きくわけて各コンピュータに得意な仕事(ベクトルコンピュータには重い計算, グラフィック・ワークステーションには途中経過の表示など)を割り当てる機能中心の並列化と, 一台のマシンでは扱えないような大きなデータを分割するデータ中心の並列化の二種類があります。今回行ったトルクコンバータの計算では, 各部品ごとに一台のワークステーションを割り当て, プログラム自体は同じというデータ並列型を用いています。またプロセス管理についても, ひとつのプロセスが他のすべてを管理するマスター・スレーブ型と, 複数の対等なプロセスが互いに協力して処理を行う SPMD (Single Program Multiple Data) があります。今回は C++ によるマスタープロセスを一つ起動し, 計算は FORTRAN によるスレーブプログラムが行う形式をとっています。

2. PVM の入手とインストール

PVM の入手方法は, anonymous ftp で持ってくるのが

*東京大学生産技術研究所 第2部

一番簡単でしょう。オリジナルは ftp://cs.utk.edu/pub/xnetlib/pvm3 にありますが, 国内では ftp.kyoto-u.ac.jp などのミラーサイトがありますので, できるだけこちらを利用するようにしましょう。それ以外にも netlib にメールを出して shar (シェルアーカイブ) 形式やエンコードされたファイルを送ってもらう方法もあります。

ftp で持ってきたファイルは gzip で圧縮してあると思いますので,

```
% gunzip pvm.3.3.8.tar.gz
```

```
% tar -xvf pvm.3.3.8.tar
```

とすれば pvm3 というディレクトリができます。この文章では pvm3 がホームディレクトリにあるものとして話をすすめていきます。まず環境変数の設定が必要になりますが, C シェルを使っていれば pvm3/lib/cshrc.stub を自分の .cshrc に追加することによって簡単に設定できます。

```
% cat pvm3/lib/cshrc.stub >> ~/.cshrc
```

```
% source ~/.cshrc
```

あとは pvm3 の下で make するだけです。

```
% cd ~/pvm3
```

```
% make
```

これで lib ディレクトリに実行ファイルができるのですが, 注意する点として, PVM を走らせるすべてのマシン上でこの make を行わなければなりません。lib ディレクトリの下には各機種ごとにサブディレクトリができますが, これらは lib/pvm などのシェルスクリプトが自動的に機種を判断して呼び出してくれますので, パスを通しておく必要はありません。

これでインストールは完了です。PVM はそのサイトのどのユーザーでもインストールすることができますが, たくさんインストールしてもディスク容量が無駄ですので, 管理者に頼んで /usr/local/pvm3 などにインストールしてもらうといいでしょう。

3. PVM による並列化の実装

次に、プログラム中で PVM をどう使用していくかというのですが、これは最初に述べた機能並列・データ並列のどちらを用いるのか、プロセスをマスター・スレーブ型にするのか SPMD にするのかによって異なります。ただ PVM はメッセージパッシングによって並列化を行いますから、データ並列を用いるほうが実装しやすいでしょう。数値計算では、計算領域や巨大行列の分割などに当たります。

まず最初にユーザーによって起動されるプログラム（マスタープログラムもしくは SPMD の最初のひとつ）が子プロセスを生成するところから始まります。デフォルトではプロセスの生成場所は PVM によって自動的に決められるのですが、現在の負荷状況を判断するといった高度なことまではしてくれませんので（初期設定ファイルに各マシンの相対的な速さを書いておくことによって少しは改善しますが）、ホスト名を指定して生成させる方が効率的です。ですから動的にプロセス数を増減させるようなもの（領域分割で分割数が多いものなど）では、各ホストの負荷状況を調べるなどの高度なプログラミングが必要となります。また生成される実行ファイルは特定の場所においておく必要があります（文献 1）参照。

データのやり取りを行う手順ですが、バッファにパックして相手プロセスに送り、受け取った側はそれをアンパックして使うということになります。

<C/C++>

```
/* バッファの初期化 */
pvm_initsend (PvmDataDefault) ;
/* 整数をパック */
pvm_pkint (&n, 1, 1) ;
/* 要素数10の実数配列をパック */
pvm_pkdouble (vec, 10, 1) ;
/* タグを付けて送信 */
pvm_send (tid, msgtag) ;
```

<FORTRAN>

```
CALL PVMFINITSEND (PVMDEFAULT, INFO)
*FORTRAN ではパックする関数はひとつだけ
CALL PVMFPACK (INTEGER4, N, 1, 1, INFO)
CALL PVMFPACK (REAL4, VEC, 10, 1, INFO)
CALL PVMFSEND (TID, MSGTAG, INFO)
```

ここで、ヴァーチャルマシンを構成するコンピュータがヘテロジニアスな場合、データをパックする際に XDR という共通のデータエンコーディングを行いますので、多少オーバーヘッドが生じてしまいます。常に同じ機種間でしか PVM を走らせないのならば、上記の PvmDataDefault の部分を PvmDataRow にすることによってそれを防ぐことができます。また、C と FORTRAN を同時に用いる場合に注意しなければならないのは、文字列の送受信です。これは C と FORTRAN で文字列の扱いが異なるために pvm_pkstr (...) と PVMFPACK (STRING,...) の内部仕様が異なる（実際 PVM の FORTRAN サブルーチンは C の関数を呼び出している）ことによるもので、数値計算ではファイル名くらいしか使わないと思いますが注意してください。

データを受け取った側では、送ったときと全く同じ順序でアンパックし、自分の計算に使用します。

<C/C++>

```
/* データがくるまで待つ */
pvm_recv (tid, msgtag) ;
/* データをアンパック */
pvm_upkint (&n, 1, 1) ;
pvm_upkdouble (vec, 10, 1) ;
```

<FORTRAN>

```
CALL PVMFRCV (TID, MSGTAG, INFO)
CALL PVMFUNPACK (INTEGER4, N, 1, 1, INFO)
CALL PVMFUNPACK (REAL4, VEC, 10, 1, INFO)
```

また、PVM 使用時にはユーザーがコンソール上で起動したプロセス以外の全プロセスの出力は、そのコンソールのあるホストのファイルにリダイレクトされます。これはこちらが意図しなくてもデータ送信が行われてしまうこととなりますので、チェックライトはデバッグ後なるべくコメントアウトする方がいいでしょう。

4. PVM をサポートするツール類

PVM は情報分野で多く用いられているため、他の並列化ツールに比べてプログラミングを支援するツール類が充実しているのも大きなメリットです。プログラマが直接行う並列化は自由度は高いのですが、デバッグの際にどのホストのどのプロセスがエラーを出しているのかわかりにくくなってしまう場合があります。また PVM のようにメッ

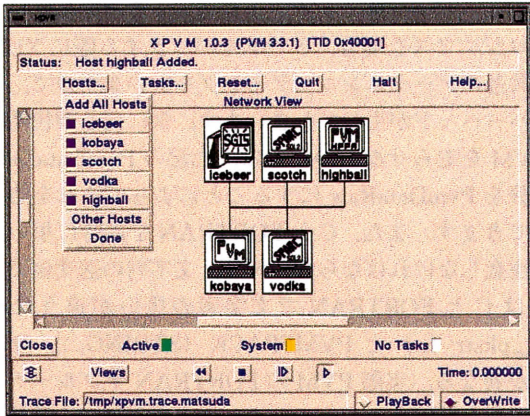


図1 X P V Mによるホストの追加

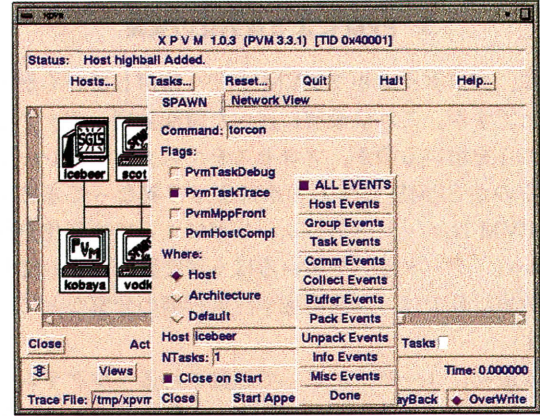


図2 X P V Mによるプロセスの始動

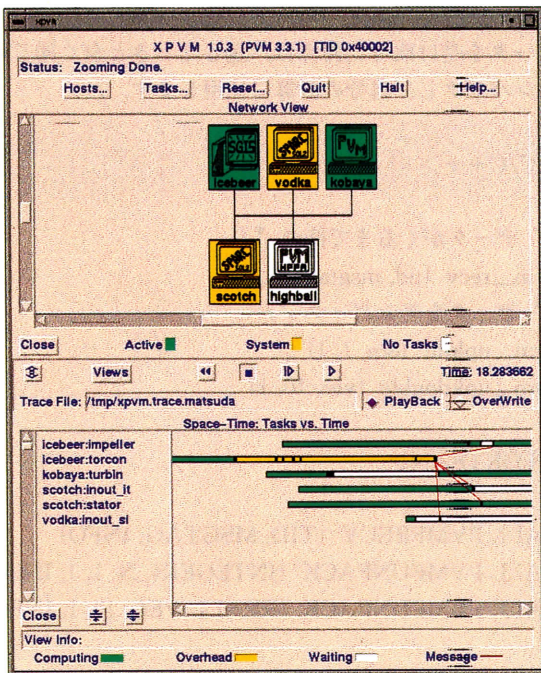


図3 ネットワークの構成と実行の様子

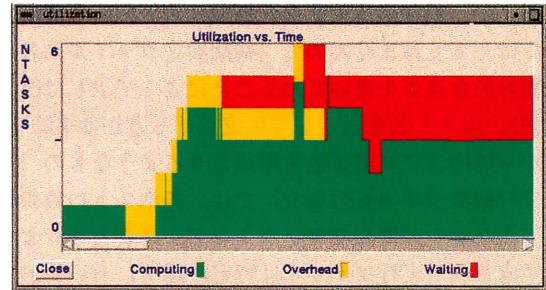


図4 プロセスの状況の表示

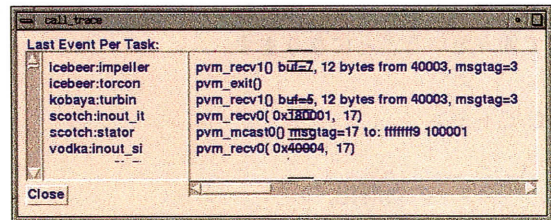


図5 コールトレース

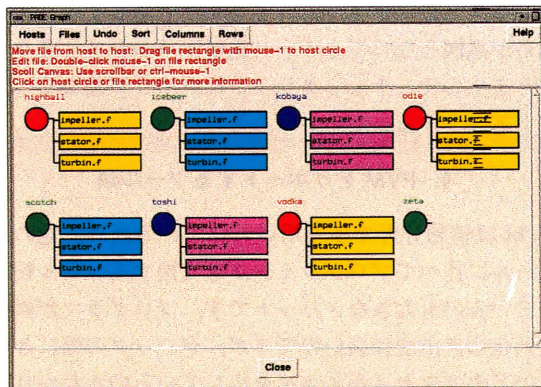


図6 PADEによるプログラム管理

セージパッシングを用いて並列化を実装しているツールでは、あるメッセージがどのプロセスからどのプロセスへ何番のタグをつけて送信されたのか、このプロセスはどのプロセスからくどのタグのメッセージを待っているのかなどをトレースできれば、デッドロックや無駄な処理待ちを防ぐことができます。またプログラミングの段階でも、多くのソースファイルを異なる各機種ごとにコピーし、デバッグが行われるごとに確実に最新のファイルを全機種でコンパイルしなおさなければならず、その手助けとなるツールもぜひ必要です。

PVMにはこのどちらもGUI(Graphical User Interface)を用いた便利なツールがありますので、デバッ

研究速報

グが容易になっています。まず実行時のプロセスやメッセージのトレースには、XPVM というツールを使用しました。これはスクリプト言語 Tcl/Tk がベースとなっており、画面上でホストの追加や削除、プロセスの実行、メッセージや各プロセスの出力のトレースなどを行うことができます。PVM 関数のコールトレースなどまでできますが、かなり重い処理になってしまうのと、並列処理という性質上各プロセスからのメッセージが非同期であることなどから、リアルタイムに処理を確認するのは困難です。そこで XPVM ではまずプロセス全体の挙動をログファイルに記録しておき、あとからコマ送りで再生することができますようになっていました。もちろん各トレースもみることができますので、デバッグに非常に便利です。

またプログラミング段階でのツールとしては、NIST-PADE²⁾ (National Institute of Standard and Technology Parallel Applications Development Environment) というパラレルメイクツールがあります。通常ヘテロジニアスな環境でのネットワークといっても、プログラミングは一台のマシンで行い、それを各機種にコピーしてメイクしなおすのが一般的です。PADE は GUI 上でそれら一連の作業を行うことができます。コピーするファイル名や実行するコマンド (make など) を登録しておけば、メニューから自動的に PVM を起動して各機種ごとにメイクを行ってくれます。PVM 自身にも aimk というメイクツールが付属していますが、PADE を用いればより簡単にプログラムを構築することができます。

コンパイラまかせの並列化とは違い、人間が直接プログラミングからデバッグまで行う必要があるわけですから、こういった開発支援ツールのよしあしが並列化ツールの普及の決め手になるといってもいいでしょう。

5. ま と め

PVM による並列化の実際について紹介してきましたが、ネットワークでつながってさえいれば大規模な計算が取り扱えるというメリットはかなり大きいと思います。プログラミングやデバッグも前出のようなツールがあることによって、まったく手の出ないものではなくなりつつあります。数値計算においては、多くのプロセッサとメモリをもつ超大型汎用並列計算機とは違う方向の、研究室規模で行える並列処理としてこれから普及していくと思われます。

(1995年11月28日受理)

参 考 文 献

- 1) Beguelin, A.L., Dongarra, J.J., Geist, G.A., Jiang, W.C., Manchek, R.J., Moore, B.K., Sunderam, V.S., Parallel Virtual Machine System, Univ. of Tennessee, Oak Ridge National Lab., Emony Univ. (1992).
- 2) Devaney, J.E., Lipman, R., Lo, M., Mitchell, W.F., Edwards, M., Clark, C.W., The Parallel Applications Development Environment (PADE) User's Manual, Computing and Applied Mathematics Lab., Physics Lab., National Institute of Standards and Technology (1995).