

Ninf による広域分散並列計算

中田 秀基^{†1} 高木 浩光^{†5} 松岡 聡^{†2}
長嶋 雲兵^{†4} 佐藤 三久^{†3} 関口 智嗣^{†1}

ローカルなネットワーク上でのメッセージパッシングライブラリを用いた分散並列計算はすでに広く行なわれている。しかし、ネットワークの高速化によって現実的になりつつある広域ネットワーク上での分散並列計算については、ソフトウェアの枠組が未だ十分に整備されていない。

我々は、広域分散並列計算に適した分散計算の枠組として「Ninf」を提案している。Ninf は広域分散環境でのマクロデータフローによる並列実行を支援するシステムで、広域での動的負荷分散とスケジューリングを特徴とする。メッセージパッシングライブラリを用いた手法に比較して (1) 広域ネットワークに適した通信パターンを用いる、(2) ユーザにとってプログラミングが容易でかつ再利用性が高い、(3) 既存のライブラリの再利用が容易、(4) ネットワーク上の資源の利用が可能、といった特長をもっている。

Global Parallel Computation using Ninf

HIDEMOTO NAKADA,^{†1} HIROMITSU TAKAGI,^{†5}
SATOSHI MATSUOKA,^{†2} MITSUHIISA SATOH^{†3}
and SATOSHI SEKIGUCHI^{†1}

Distributed computing using message passing libraries in a LAN(Local Area Network) environment is already accepted as an effective supercomputing methodology. On the other hand, although distributed computing in WAN(Wide Area Network) environment is becoming practical due to recent development of high-speed network facilities, software framework for supercomputing in WAN is yet to be established. We propose 'Ninf', a distributed computing framework

for globally distributed computing environment. Ninf enables parallel computing in WAN based on the macro dataflow model, and facilitates automatic dynamic load distribution and scheduling. Ninf has the following advantages over using existing message passing libraries in WAN supercomputing: (1) communication protocol suited for globally distributed environment, (2) ease of programming (3) reuse of existing libraries, (4) integration with existing data resources on the Internet.

1. はじめに

科学技術計算に要求される計算の規模は増大する一方である。これに対して、計算機の単体性能は向上しているとはいえ、使用電力や計算機の構成にかかる金銭的な問題から、単一サイトにおける計算機の能力は頭打ちになると考えられる。したがって、要求される計算の規模を満足するだけの計算能力を提供するためには、高速なネットワークを用いて複数のサイトに分散して存在する計算機に計算を分散することが必要になる。現在、LAN 内のワークステーションを用いた分散計算はすでに広く行われているが、広域分散環境では通信速度と計

算速度の比率が LAN 内とは著しく異なるため別の枠組が必要となる。

我々は、広域分散並列計算に適した分散計算の枠組として「Ninf」を提案している。Ninf は広域分散環境でのマクロデータフローによる並列実行を支援するシステムで、広域での動的負荷分散とスケジューリングを特徴とする。メッセージパッシングライブラリを用いた手法に比較して (1) 広域ネットワークに適した通信パターンを用いる、(2) ユーザにとってプログラミングが容易でかつ再利用性が高い、(3) 既存のライブラリの再利用が容易、(4) ネットワーク上の資源の利用が可能、といった特長をもっている。

本稿では、Ninf におけるユーザアプリケーションのプログラミングモデルについて述べ、メッセージパッシングライブラリを利用する場合に比較して、容易性、再利用性、汎用性、拡張性が高いことを示す。また、広域ネットワーク上での分散計算において効率的な計算を実

†1 電子技術総合研究所 Electrotechnical Laboratory

†2 東京工業大学 Tokyo Institute of Technology

†3 新情報処理開発機構 Real World Computing Partnership

†4 お茶の水女子大学 Ochanomizu University

†5 名古屋工業大学 Nagoya Institute of Technology

現できることを予備的な評価によって示す。

2. 広域分散並列計算の要件

ワークステーションを LAN で結合したワークステーションクラスタによる分散並列計算は現在広く行なわれている。これらの多くは、PVM や MPI などのメッセージパッシングパラダイムに基づくシステムを用いている。このパラダイムによる記述には、プログラム間のデータのやりとりの自由度が高いという利点がある反面、データ交換のプロトコル定義やメッセージの明示的な送受信の記述をユーザーが行う必要があるという欠点がある。

広域分散環境下での分散並列処理が、LAN 上の並列処理と比較して最も大きく異なるのは、単一計算機の計算速度に対するネットワークを介した通信の速度の比率が著しく低いという点である。広域ネットワークの通信速度は向上しつつあるとはいえ、レイテンシに関しては劇的な向上は見込めない。一方で、計算速度は今後も向上し続けることが予想される。このように計算速度が通信速度に比較して非常に高速な場合には、通信のコストが相対的に著しく高くなり、効率的に並列実行できるアプリケーションの型が限定される。短い期間ごとに同期をとりデータを交換する、例えば LU 分解のようなアプリケーションを広域分散環境下で高速化することはむずかしい。計算要素、通信の粒度がともに大きいアプリケーションでなければ効率的な実行は望めない。

科学計算の分野には、このような粗粒度での並列実行に適したアプリケーションが数多く存在する。特定のパラメータのある範囲で走査しながらそれぞれの地点での値を計算するような問題、非常に大量の乱数値に対して特定の演算を行って確率的な解を得るモンテカルロシミュレーションのような問題などである。これらの問題は、多数の相互通信を必要としない独立した計算で構成されており、基本的にオーバーヘッドなく任意のサイズに分割し並列に実行することができる。

このような問題を粗粒度並列で記述するには、メッセージパッシングのようなパラダイムは不適である。メッセージパッシングは、データ交換の頻度がある程度高く、複雑な通信形態が必要とされる場合を前提にしている。このため、非常に簡単な通信形態を記述する場合にも、プロトコルの設計や明示的なメッセージ送受信をユーザがする必要がある。これは、ユーザにとって大きな負担となる。簡単な通信形態を、簡単に記述できる枠組みが必要である。

従来のメッセージパッシングシステムのもう一つの問題点として、構成要素計算機と通信系の信頼性の問題がある。従来のメッセージパッシングシステムは、すべての計算機が同一 LAN 内に存在し、それらすべてがある程度信頼できる形で運用されていることを前提としている。しかし広域分散環境では、通信回線の信頼性は低

く、各々の計算機の管理体制も異なることを前提としなければならない。広域分散環境でのシステムには、このような環境でも頑健に運用できることが要求される。PVM や MPI などのメッセージパッシングライブラリでは、通信相手となるホストを適切に広域環境上で指定する方法は定義されておらず、かつ対故障性なども考慮されていないので、これらの要求は満たされていない

3. Ninf による広域分散並列計算

3.1 Ninf システムの概要

我々は、前節で述べたような要求を満たすシステムとして、Ninf システムを設計した¹⁾。Ninf は計算を粗粒度な要素ルーチンに分割してネットワーク上に分散した計算機に配置し、それらの間にデータフローを構成し、計算を行なうシステムである。

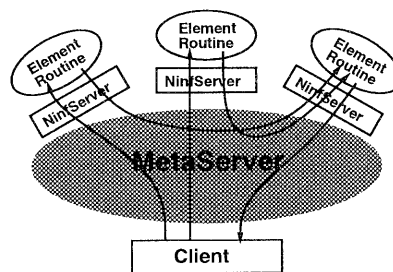


図1 Ninf システム

Ninf システムは、クライアント、Ninf サーバ、要素ルーチン、メタサーバの4つの要素からなる。要素ルーチンは独立したプログラム部品で、再利用可能な一般的なインターフェイスを与えて設計する。要素ルーチンは複数の入力を得、複数の出力を与えるだけのもので、状態は持たない。状態を持たないため、要素ルーチンのインターフェイスは非常に単純になる。この要素ルーチンを組み立てる枠組みとして、クライアントプログラムを記述する。クライアントプログラムには、使用する要素ルーチンとその間のデータ依存関係を記述する。データ依存関係が最も単純な場合、すなわちクライアントと一つの要素ルーチンの間にしか依存関係がない場合が通常のクライアント / サーバモデルになる。

Ninf サーバは分散システムを構成する各々の計算機上で待機し、それぞれの計算機上の要素ルーチンを管理する。リクエストに応じて個々の要素ルーチンを実行する役割を持つ。メタサーバは、クライアントと Ninf サーバを繋ぐ役割を果たす²⁾。クライアントプログラムに記述された要素ルーチン間のデータ依存関係は、クライアントの Ninf ライブラリによってメタサーバに転送される。メタサーバはデータ依存関係に基づいて、スケジューリング、負荷分散を行ない、Ninf サーバへリク

エストを発行し要素ルーチンを実行する。要素ルーチン間のデータ転送もメタサーバの役割の一つである。

各計算要素ルーチンの記述は、通常の C または Fortran の関数の記述と同様に行なうことができる。既存のライブラリに簡単なインターフェイス記述を加えるだけでそのまま使用することができるので、各専門領域における膨大なソフトウェア資産を活かすことができる。各要素ルーチン間には、要素ルーチン同士がメモリ空間を共有しているかのように見せるプログラミングインターフェイスが提供される。データフローの上流の計算要素が計算結果を自分のメモリに書き込むと、その値が下流の計算要素のメモリに自動的に転送される。このため、メッセージパッシングのように、プログラマが明示的に通信を記述する必要はなく、各計算要素は通常の単一ノードでの逐次プログラムと全く同じように記述することができる。このような簡潔なインターフェイスにすることによって、プログラマの負担が軽減されている。通信プロトコルの設計ミスによるデッドロックなどのバグが発生する可能性はない。

Ninf では問題の記述を、要素ルーチン自身の記述とデータ依存関係の記述という 2 つのレベルに明示的に分離する。このため、システムとしての汎用性、拡張性に富んでいる。例えば、要素ルーチンの各入力ソースは、プロトコルを満たしてさえいれば何でもかまわない。たとえば、WWW や ftp から取得したデータを、計算の入力として直接利用することも可能である。

3.2 インターフェイスの記述

要素ルーチンのインターフェイスは Ninf IDL (Interface Description Language) と呼ばれる言語で記述する。この記述をインターフェイスジェネレータでコンパイルして、Ninf の提供するライブラリとのインターフェイスの役割を果たすスタブルーチンを作成する。このスタブと関数本体をリンクして、要素ルーチンの実行形式を得る。Ninf サーバはリクエストに応じてこの実行形式を実行する。

3.3 Ninf のプログラムインターフェイス

Ninf クライアントは現在、C、Java、Lisp 環境に実装されているが、ここでは C 言語の環境を説明する。

Ninf システムを用いて計算をリモートの計算機で実行するには、以下のよう書く。

```
Ninf_call(func_name, .....);
```

Ninf_call の第一引数には URL に準じた記法でサーバのホスト名 / ポート、計算を指定する。

```
ninf://hostname[:port]/funcname
```

サーバ、ポート名を省略し、関数名だけを記述した場合にはシステムのデフォルトとして指定されたメタサーバが、適切なホストを選択する。

第 2 引数以降には通常の C の関数呼び出しと同様に引数を記述する。

この Ninf 呼び出しはブロッキングである。すなわち、リモート計算機での計算が完全に終了し、計算結

果を回収するまでリターンしない。これでは、複数の Ninf 呼び出しを同時に発行することができないので、計算の発行と計算結果の回収を分割して行なう方法が用意されている。

```
Ninf_call_async(func_name, ...);
```

この関数は、要素ルーチンの実行と引数の送信を行ない、計算の終了を待たずにリターンする。返り値として、その Ninf 呼び出しのセッションの ID を返す。

計算の結果を回収するには、この ID を引数として

```
Ninf_wait(ID);
```

を呼び出す。このルーチンはそのセッションの計算が完了するのを待ち、すべての結果を回収し Ninf_call_async の引数に書き込んでリターンする。これらを用いると、

```
id1 = Ninf_call_async(func_name, ...);
```

```
id2 = Ninf_call_async(func_name, ...);
```

```
;
```

```
Ninf_wait(id1);
```

```
Ninf_wait(id2);
```

```
;
```

のように記述して、複数の計算を並行して発行することができる。

すべてのセッションの終了を待つこともできる。

```
Ninf_wait_all();
```

とすると、計算が終了したセッションから順に結果を回収し、すべてのセッションが終了するまで待つ。

3.3.1 トランザクション

Ninf_call_async を用いれば互いに依存関係のない単純な並行計算は容易に記述できる。しかし Ninf 呼び出しの間に依存関係がある場合には、適切な場所に同期を挿入することが必要となる。これは並列プログラミングに不馴れなプログラマにとっては容易ではない。そこで、データ依存関係にしたがった発行と同期を自動化するためにいくつかの連続する Ninf 計算呼び出しをまとめた、トランザクションという機構を導入した。トランザクションは開始点と終了点を指定することで、ユーザによって明示的に設定される。トランザクション内の複数の Ninf_call は依存関係が自動的に解析され、その依存関係が保証されるようにスケジュールされ、並列実行される。

例えば、行列 A、B、C、D の加算は以下のように単純に記述できる。

```
Ninf_transaction_begin();
```

```
Ninf_call("madd", n, A, B, E);
```

```
Ninf_call("madd", n, C, D, F);
```

```
Ninf_call("madd", n, E, F, G);
```

```
Ninf_transaction_end();
```

実際の計算の発行は、トランザクション終了時に一括して行なわれる。トランザクションの内部では、Ninf_call のセマンティクスが変化し、実際の Ninf 呼び出しを行わず、データ依存グラフの構築のみ

を行い、すぐにリターンするようになる。

トランザクションが終了すると、クライアントのライブラリが、データ依存グラフと各 Ninf 呼び出しにおける引数とをメタサーバに転送する。メタサーバは、データ依存グラフに従って、各 Ninf 呼び出しをスケジューリングし、Ninf サーバへのリクエストを行なう。この例の場合 1 番め、2 番めの Ninf 呼び出しが終了したのち、3 番めの Ninf 呼び出しが発行される。

C 言語環境でのデータ依存解析は、配列データのポインタの同一性と Ninf 呼び出しの各引数の入出力モードに基づいて行なう。前述したように、Ninf の各要素ルーチンは自分のインターフェイス情報を持っており、Ninf 呼び出し時にこの情報が動的に転送される。このインターフェイス情報にモード情報も含まれており、データ依存解析にはこの情報を用いる。例の場合、madd の第 3、4、5 引数の入出力モードはそれぞれ、入力、入力、出力であり、図 2 のようなデータ依存関係があることがわかる。

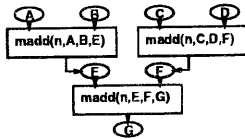


図2 行列加算のデータ依存関係

このとき中間データのクライアントへの転送は抑制される。上記の例でいえば、変数 E, F が中間データである。これらの値はデータ依存関係を表現するために用いられたに過ぎず、クライアントプログラムで必要とされる値ではないためクライアントに書き戻さない。

3.4 メッセージバッシングライブラリとの比較

Ninf による記述と代表的なメッセージバッシングライブラリの一つである PVM による記述とを比較してみよう。ここでは、NAS Parallel Benchmark の一つで、モンテカルロ法によって円周率を求める計算である EP を記述してみる。

図 3 に Ninf による EP のコア部分を示す。

```

Ninf_transaction_begin();
for (i = 0; i < PU; i++){
    Ninf_call(buffer, i, NPP, ktmp+i,
             xsutmp+i, ysutmp+i, qtmp[i]);
}
Ninf_transaction_end();

```

図3 Ninf による EP の実行

ループの手前でトランザクションを開始し、ループの中で複数の Ninf 呼び出しを行い、ループの後でトランザクションを終了している。Ninf では、このように簡潔なプログラムを記述するだけで、ep のサブプログラ

ムが自動的に適切な要素計算機上に分散され、並列に実行される。

これに対して、PVM ではプログラマが、プロセスをフォークするコード、解を回収するコードをいちいち手で記述しなければならない。図 4 に PVM でのスレーブプロセスのフォークのコードをしめす。

```

for(i=0; i<PU-1; i++){
    if(pvm_spawn(PROG,NULL,TASK_CREATE,
                 hname[(i%hnum)],1,&tids[i+1])!=1){
        perror("Cannot create process\n");
        pvm_exit();
    }
}

```

図4 PVM でのフォーク

フォークの際にはホストネームを明示的に与える必要がある。さらに、ここでは省略するが解の回収のコードも必要となり、本来プログラムの本質的でない部分である通信部分の記述がプログラマにとって大きな負担になることがわかる。

もう一つの問題点は、再利用性である。メッセージバッシングライブラリでは、プログラムの用いているプロトコルはプログラムに非明示的に埋め込まれてしまうため、要素プログラム自体に一般性があったとしてもそれを再利用することが難しくなる。Ninf では、要素プログラムのインターフェイスは IDL ファイルとして独立して記述されるため、要素プログラムの再利用は非常に容易である。

4. 評価

Ninf による分散並列実行の予備評価として、LAN 内での並列実行を評価した。

4.1 実験環境

実験環境としては、DEC Alpha (333 MHz) の 32 台のワークステーションクラスタを用いた。各ワークステーションは、100Base/T のスイッチで結合されている。32 台の他にクラスタのホストとして 1 台の Alpha ワークステーションがある。これらは NFS でファイルを共有している。

Ninf サーバをそれぞれのワークステーションで実行し、クライアントには、ホストのワークステーションを用いた。今回はメタサーバを用いず、非同期の Ninf 呼び出しを用いて並列実行した。

4.2 EP による評価

前述した EP を用いて評価を行った。この計算では、乱数による試行は完全に独立に行なうことができ、さらにそれぞれの試行一回あたりのコストは一定であるから、容易に負荷を等分することができる性質のよいプログラムである。

実験は、SampleとClassAとClassBとで行った。Sampleでは 2^{24} 回、ClassAでは 2^{28} 回、ClassBでは 2^{30} 回の試行を行う。問題の分割法としては、単純に試行回数をプロセッサ数で割って各プロセッサに割り付けた。32プロセッサの場合各要素ルーチンでの指向回数はSample、ClassA、ClassBでそれぞれ、 2^{19} 回、 2^{23} 回、 2^{25} 回の実行が行われることになる。図5に実行結果を示す。

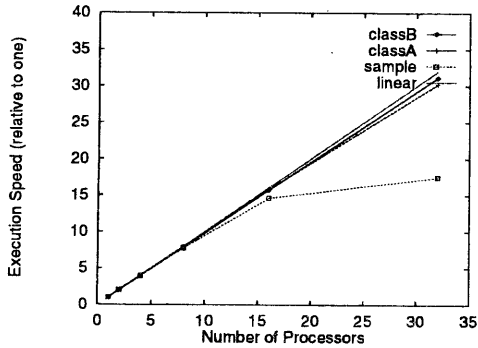


図5 EPの実行

Sampleでは、32台の時に大きく性能が低下している。ClassA、ClassBではほぼリニアに速度が向上しているが、ClassAでは32台でわずかに低下が見られる。これは問題が細分化され、個々の要素ルーチンでの実行時間が短くなった結果、Ninfによるオーバーヘッドが顕在化したものであると思われる。

分散台数を常に32台にし、個々の要素ルーチンでの試行回数を 2^{15} から 2^{25} まで変化させた結果を図6に示す。値はローカルに同じ試行回数を実行した場合との比較である。

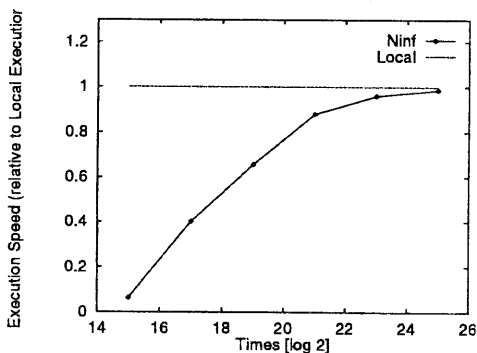


図6 試行回数と相対実行速度

Ninfの呼び出しによるオーバーヘッドは分散台数にのみ依存し、この場合は固定であると考えて良い。試行回数が増え要素ルーチンの実行時間が延びるにしたがっ

て、Ninf呼び出しのオーバーヘッドが無視できるようになるのがわかる。

4.3 状態密度計算

この計算は巨大分子の状態密度を、行列対角化することなく求める手法である。この状態密度計算 (density of states)³⁾ は、巨大行列をバネのつながった力学モデルに見立て、周期的な外力をあたえ、バネの共鳴の度合いを状態数とするものである。外力の周波数に対して並列化を行うことができるため、並列計算に適している。

図7に実行結果を示す。

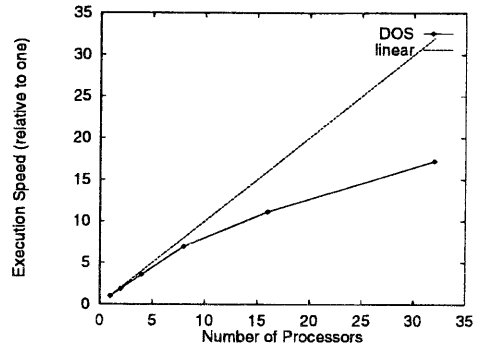


図7 状態密度計算の台数効果

問題の分割は、外力の周波数に対してサイクリックに行った。台数効果がそれほどでないのは、負荷分割が適切でなかったためであると思われる。この計算のそれぞれの周期に対する計算量は、周期に依存して大きく変化する。これを単純にサイクリックに分割してしまっているため、負荷に偏りが生じ、台数効果が上がらなかったためであると推測される。

4.4 PVM との比較

図8にPVMで記述した場合に対する、Ninfで記述した場合の実行速度の比較を示す。対象プログラムは前出のEPである。

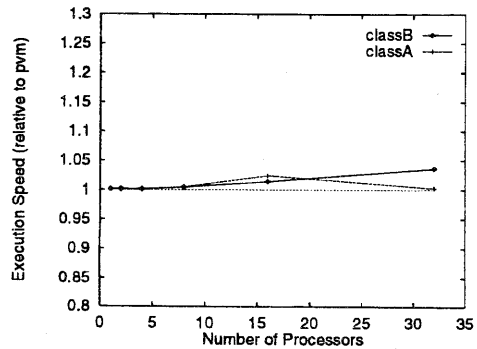


図8 PVMを用いた場合との比較

PVMによるプログラムは、問題の分割手法、通信のプロトコルがNinfによるものと同じになるよう記述した。Ninfで記述したものは、PVMで記述したものとほぼ等価か、より少ないオーバーヘッドで実行できているのが分る。

4.5 広域環境下での性能低下予測

今回の評価では、クライアントのマシンとNinfサーバのマシンとは高速なLANで接続されている。本節では、広域環境で通信速度が低下しNinfによるオーバーヘッドが大きくなった場合の実行速度について考察する。

通信速度の低下によりオーバーヘッドは増大する。しかし、問題の計算時間が充分長い場合実行速度の面ではほとんど問題にならない。先ほどのepを例にとる。ClassBを32台で演算した場合の試行回数は 2^{25} である。これを実行するには、およそ、67.3秒ほどかかる。この場合のNinfのオーバーヘッドは1秒弱に過ぎない。したがって、このオーバーヘッドが仮に10倍になったとしても、オーバーヘッドは全実行時間の1/8程度に過ぎない。

オーバーヘッドが2倍、5倍、10倍、20倍であったとしたときの、ローカルな場合との相対的执行速度を図9に示す。

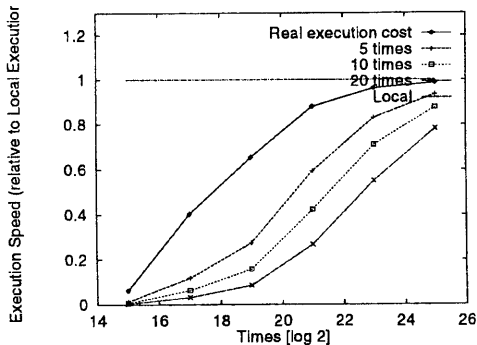


図9 オーバヘッドが大きい場合の相対実行速度

ここで示したアプリケーションは、Ninfのターゲットとしては比較的要素ルーチンの実行時間の短いアプリケーションである。しかし、このような場合でさえオーバーヘッドが増大したとしても計算時間のほうが支配的である。アプリケーションの規模が増大すれば、さらに計算時間が支配的になるので、Ninfのオーバーヘッドは広域環境下でも問題にならないと思われる。

5. 関連研究

分散環境でのスーパーコンピューティングを行うプロジェクトは、幾つか行われつつある。NetSolve⁴⁾は、Ninfによく似たシステムで、クライアントのAPIも

類似している。NetSolveにはデータフローを解析し、スケジュールする機能がない点がNinfと異なる。APIのセマンティクスはほぼ等価なので、NetSolveで用意された数値演算ライブラリをNinfから使用し、Ninfで用意した数値演算ライブラリをNetSolveから使用する、「相互乗り入れ」を検討中である。

スイスETHのRCS(Remote Computation System)⁵⁾は、Ninfと同様にバックエンドのスーパーコンピュータをリモートから使用するシステムである。Ninfとの相違としては、RCSでは下位のレイヤとしてPVMを用いており広域分散に適していないこと、動的なサーバの追加ができないこと、並列計算に対応していないことが挙げられる。また、Ninfではすべてのリモートルーチンを一つのAPIで呼び出すことができ、ルーチンの追加時にクライアントの変更が必要ないのに対して、RCSではリモートルーチンを追加することに別の関数をインターフェイスとして追加する。このため、システムの拡張性に乏しい。

6. おわりに

広域分散並列計算に適した分散計算の枠組への要請について議論し、我々のシステムNinfがその要請を満たすことを示した。また、ワークステーションクラスタを用いた予備的な評価を行い、Ninfが分散環境において効率的な計算を実現でき、PVMに遜色ない実行性能であることを示した。

今後の課題としてはより洗練された動的な負荷分散が挙げられる。

参考文献

- 1) Sekiguchi, S., Sato, M., Nakada, H. and Nagashima, U.: - Ninf - : Network base information library for globally high performance computing., *Proceedings of Parallel Object-Oriented Methods and Applications (POOMA)* (1996).
- 2) 中田秀基, 草野貴之, 松岡聡, 佐藤三久, 関口智嗣: ネットワーク数値ライブラリ Ninf におけるメタサーバアーキテクチャ, 情報処理学会研究報告 HPC, Vol. 96, No. 22 (1996).
- 3) 日向寺祥子, 長島雲兵, 細矢治夫, 関口智嗣, 佐藤三久: 大規模実対称行列の状態密度の計算とその並列化, 情報処理学会研究報告 HPC, Vol. 93, No. 93-HPC-48, pp. 57-64 (1993).
- 4) Casanova, H. and Dongara, J.: NetSolve: A Network Server for Solving Computational Science Problems, *Proceedings of Super Computing '96* (1996).
- 5) Arbenz, P., Gander, W. and Oettli, M.: The Remote Computation System, Technical Report 245, ETH (1996).