

JAERI-Data/Code
98-030



Paragon上でのスカラー超並列プログラム開発ガイド

1998年10月

上島 豊・荒川拓也・佐々木明・横田 恒*

日本原子力研究所
Japan Atomic Energy Research Institute

本レポートは、日本原子力研究所が不定期に公刊している研究報告書です。
入手の問合せは、日本原子力研究所研究情報部研究情報課（〒319-1195 茨城県那珂郡東海村）あて、お申し越しください。なお、このほかに財団法人原子力弘済会資料センター（〒319-1195 茨城県那珂郡東海村日本原子力研究所内）で複写による実費領布をおこなっております。

This report is issued irregularly.

Inquiries about availability of the reports should be addressed to Research Information Division, Department of Intellectual Resources, Japan Atomic Energy Research Institute, Tokai-mura, Naka-gun, Ibaraki-ken 319-1195, Japan.

© Japan Atomic Energy Research Institute, 1998

編集兼発行 日本原子力研究所

Paragon 上でのスカラー超並列プログラム開発ガイド

日本原子力研究所関西研究所光量子科学センター

上島 豊*・荒川 拓也・佐々木 明・横田 恒*

(1998年9月11日受理)

日本で100並列を越える並列計算が、実際に行われるようになったのは、つい数年ほど前からである。日本原子力研究所（原研）のIntel製75MP834<関西研究所>、Paragon XP/S 15GP256 <那珂研究所>は、本格的超並列計算機の先駆けとして光量子、核融合の大規模、超並列計算を目的に導入されている。これらの計算機を使って超並列計算を行うために、多くの超並列計算プログラムが移植や新規作成されている。これらのプログラムは、従来、ワークステーションやベクトル計算機上で動作していたものをそのまま移植したものか、並列用にアルゴリズムの変更を施したものである。異なる計算機及びオペレーティングシステム（OS）の基でのプログラム開発には、細心の注意とノウハウが必要であるが、Paragonに到ってはユーザ数が極めて少ないため、ノウハウの集積と環境の標準化が大変困難な状況にある。そのため、原研関西研究所におけるParagon XP/S 75MP834 上での超並列計算プログラム開発において得た情報をParagon上での超並列プログラム開発ガイドとしてまとめた。

Guide to Development of a Scalar Massive Parallel Programming on Paragon

Yutaka UESHIMA*, Takuya ARAKAWA, Akira SASAKI
and Hisasi YOKOTA*

Advanced Photon Research Center
Kansai Research Establishment
Japan Atomic Energy Research Institute
Miiminami-machi, Neyagawa-Shi, Osaka

(Received September 11, 1998)

Parallel calculations using more than hundred computers had begun in Japan only several years ago. The Intel Paragon XP/S 15GP256 <Naka Fusion Research Establishment>, 75MP834<Kansai Research Establishment> were introduced as pioneers in Japan Atomic Energy Research Institute (JAERI) to pursue massive parallel simulations for advanced photon and fusion researches. Recently, large number of parallel programs have been transplanted or newly produced to perform the parallel calculations with those computers. However, these programs are developed based on software technologies for conventional super computer, therefore they sometimes cause troubles in the massive parallel computing. In principle, when programs are developed under different computer and operating system (OS), prudent directions and knowledge are needed. However, integration of knowledge and standardization of environment are quite difficult because number of Paragon system and Paragon's users are very small in Japan. Therefore, we summarized information which was got through the process of development of a massive parallel program in the Paragon XP/S 75MP834.

Keywords: A Massive Parallel, Paragon, Guide, Standardization, Integration of Knowledge

* Post-doctoral Fellow

* Research Organization for Information Science & Technology

目 次

1. はじめに	1
2. paragon の特徴	2
3. 並列計算プログラムの特徴	4
4. paragon 上での超並列プログラム開発の指針	6
5. 結 び	25
謝 辞	25
参考文献	25
付録1 サンプルプログラム メインルーチン	26
付録2 サンプルプログラム サブルーチン	42
付録3 メッセージパッシングの手法	52
付録4 並列入出力の方法	60
付録5 割り込み処理によるプログラム実行の制御	67

Contents

1. Introduction	1
2. Character of Paragon	2
3. Feature of Parallel Programming	4
4. Guide to Massive Parallel Programming on Paragon	6
5. Conclusion	25
Acknowledgment	25
References	25
Appendix.1 Sample Programs (Main Routine)	26
Appendix.2 Sample Programs (Subroutine)	42
Appendix.3 Method of Message Passing	52
Appendix.4 How to Use Parallel Input/Output	60
Appendix.5 Control of Program Execution Using Interrupt	67

This is a blank page.

1. はじめに

ここ4、5年前までは、最も高速な計算機といえば、いわゆるスーパーコンピュータと呼ばれるベクトル型計算機であった。この型の計算機は4、5年で計算速度が10倍ずつ高速になってきた。しかし、現在、1CPUの演算速度が限界に達してきている。今まで使われていたベクトル計算機の1CPU性能は、数GFLOPSで搭載メモリも数Gbyteであったが、将来、格段に速い単体の計算機は出てきそうにない。これにかわって、より高速な計算機として並列計算機が次々と導入されてきている。

日本原子力研究所（原研）には、Intel製 Paragon XP/S 75MP834 <関西研究所^{1, 2)}、15GP256 <那珂研究所>が、本格的超並列計算機の先駆けとして光量子、核融合の大規模、超並列計算を目的に導入されている。これらの計算機を使って超並列計算を行うために、多くのプログラムの移植や新規作成が行われている。これらのプログラムは、従来、ワークステーションやベクトル計算機上で動作していたプログラムをそのまま移植したものか、並列用に多少のアルゴリズム変更を施したものである。一般的に、異なる計算機及びオペレーティングシステム（os）の基でプログラムを開発するには、細心の注意とノウハウが必要であるが、従来、ワークステーションからベクトル計算機、異なるメーカーの計算機へのプログラムの移植は、ユーザ数の多さが幸いして、多くのノウハウと標準化された環境のおかげで、比較的短時間で作業を完了することができていた。しかし、超並列計算機に関しては、ユーザ数が極めて少ないため、ノウハウの集積と環境の標準化が大変困難な状況にある。そのため、原研関西研究所における Paragon XP/S 75MP834 上での超並列計算プログラム開発において得た情報を Paragon 上での超並列プログラム開発の指針としてまとめた。

関西研究所で稼働中の分散メモリシステム超並列計算機 Paragon XP/S 75MP834 は、システム全体で 120GFLOPS、100Gbyte の性能を有している。この数字は、現有の単体の最高速計算機の約 100 倍の演算速度と搭載メモリの性能を意味するものである^{3, 4)}。実際、この性能がどの程度のものかを理解してもらうために、シミュレーションしている系（レーザー・プラズマ相互作用）で例を挙げる。従来の単体ベクトル計算機では、典型的なシミュレーションで粒子数が数千万粒子、空間格子が $1,000 \times 1,000$ 格子のシミュレーションが限界であった。しかし、この Paragon 上で並列化されたプログラムを使うと、粒子数が数十億粒子、空間格子が $1,000 \times 100,000$ 格子または $10,000 \times 10,000$ 格子のシミュレーションが実行できる⁴⁾。このシステムサイズの増大がもたらす波及効果は、絶大なものである。もし、空間格子を $0.1\mu\text{m}$ （レーザー波長の $1/10$ ）とすると、 $100\mu\text{m}$ （レーザー広がり方向） $\times 10\text{cm}$ （レーザー伝搬方向）という、実際の実験に近いレーザーの伝播を計算することができる。また、2次元計算が困難とされていた固体（電子密度 $\sim 10^{23}/\text{cm}^3$ ）との相互作用では、初期温度 100eV で $5\mu\text{m}$ （レーザー広がり方向） $\times 1\mu\text{m}$ （固体の厚さ）のレーザーの反射（高調波発生）吸収のシミュレーションも行うことができる。

2. Paragon の特徴

Paragon とは、Intel 社製の i860XP プロセッサを使った分散メモリスカラー超並列計算機群の総称名であり、原研では那珂研究所（那珂研）と関西研究所（関西研）に、それぞれ Paragon XP/S 15GP256、Paragon XP/S 75MP834 が稼働中である。

プログラム開発に関して関西研の Paragon を使用したため、以下、関西研の Paragon XP/S 75MP834 についてさらに詳しく述べておく。関西研の Paragon XP/S 75MP834 のハードウェア仕様は、Table.3.1 の通りである。

Table.3.1 Pragon XP/S-75MP834 と分割後の S120MP、S5MP のハードウェア仕様

	XP/S-75MP834	S120MP	S5MP
Compute ノード名	MP ノード	MP ノード	MP ノード
Compute ノード数	834	800	34
1 ノードの演算プロセッサ数	2×i860XP	2×i860XP	2×i860XP
演算性能 1 ノード	150MFLOPS	150MFLOPS	150MFLOPS
全 ノード	125.1GFLOPS	120GFLOPS	5.1GFLOPS
Memory 1 ノード	128MB	128MB	128MB
全 ノード	106.752GB	102.4GB	4.352GB
ノード内転送速度	400MB/s	400MB/s	400MB/s
ノード間結合方式	2 次元 Mesh	2 次元 Mesh	2 次元 Mesh
ノード間転送速度	175MB/s	175MB/s	175MB/s
磁気ディスク	432GB 20GB(ufs) 240GB(pfs)	10GB(ufs) 8GB(pfs)	

この Paragon XP/S 75MP834 は、バッチジョブ・大規模計算用の Paragon S120MP (800 ノード) とインタラクティブ・バック・小規模計算用の Paragon S5MP (34 ノード) にハードウェア的にもソフトウェア的にも分割され、ともに独立に運用されている。1 MP ノードの演算性能は、150MFLOPS (= 75 MFLOPS / i860XP × 2) であり、メモリ容量は 128Mbyte である。システム全体では、Paragon S120MP は、120GFLOPS、102.4Gbyte、また、Paragon S5MP は、5.1GFLOPS、4.352Gbyte である。また、MP ノード内のメモリと i860 XP プロセッサ間のデータ転送速度は、400MB/s である。MP ノード間は、2 次元相互接続ネットワークアーキテクチャによって結合しており、隣り合った任意の MP ノード間のデータ転送速度は、送受信とも 175 MB/s、双方向で 350MB/s である。

磁気ディスク装置は RAID3 を採用し、Paragon S120MP の磁気ディスク容量が 400Gbyte、Paragon S5MP の磁気ディスク容量が 32Gbyte である。また、この磁気ディ

スクは、通常の Unix File System (ufs=work) 領域の他に、並列計算機で発生する大容量データの入出力を高速に行うための Parallel File System (pfs) という並列 I/O が可能な領域を持っている。Paragon から実際に利用できる磁気ディスクシステムは、上記に述べた RAID3 であり Paragon S120MP に ufs 領域 20Gbyte と pfs 領域 240 Gbyte、Paragon S5MP に ufs 領域 10Gbyte と pfs 領域 8Gbyte に分割されている。その他に Network File System により接続されたファイルサーバ (S-4/1000E) があり、RAID5 フォーマットの ufs 領域 210G byte を利用できる。

搭載されているソフトウェアは、Table.3.2 の通りであり、2つの Paragon は、UNIX に準じた Paragon 用マルチコンピュータ・オペレーティングシステム OSF/1 により制御されている。

Table.3.2 Pragon ソフトウェア仕様

ソフトウェア名	概要
Paragon OSF/1	Paragon 用 UNIX OS
if77	Paragon 用 Fortran 77
icc	Paragon 用 ANSI C
iCC	Paragon 用 C++
ihpf	Paragon 用 HPF
NX ライブラリ	並列化のための基本ライブラリ
BLAS 1	数値計算ライブラリ Vector-Vector
BLAS 2	数値計算ライブラリ Matrix-Vector
BLAS 3	数値計算ライブラリ Matrix-Matrix
Pro Solver Parallel FFT	数値計算ライブラリ FFT
IPD	並列化プログラムデバッカ
XIPD	x 対応の並列化プログラムデバッカ
Performance Analyzer	プログラム並列化支援性能評価ツール

また、並列計算機ではあるが、個々の計算ノードに UNIX OS が動いているのではなく、個々の計算ノードには、分散 OS (Mach3.0 マイクロカーネル) が搭載され、分散 OS の UNIX に透過的なインターフェイスを介して、サービスノードに搭載されている UNIX OS の機能を利用する構造となっている。このサービスノードは、全計算用ノードの制御も行っている。また、並列化に必要なライブラリとして NX ライブラリが標準でサポートされており、この NX ライブラリをベースにしたメッセージパッシングライブラリとして MPI が使用できる。

3 並列計算プログラムの特徴

超並列計算機へのプログラムの移植は、まだ始められて間もないため、経験の集積と環境の標準化がほとんどないが、構造化プログラミング、イベント駆動プログラムなどの最近のプログラム技術を活用することで、多くの問題を解決できると考えられる。

並列化の場合は、従来のシングルプロセッサの計算機とは異なり、プロセッサ間データの共有、同期や I/O アクセスについて配慮が必要である。Paragon の構成は Fig.4.1 に示すとおりであり、計算ノードが相互にメッセージパッシングの方法で、データの授受および I/O を行いながらプログラムを実行する。

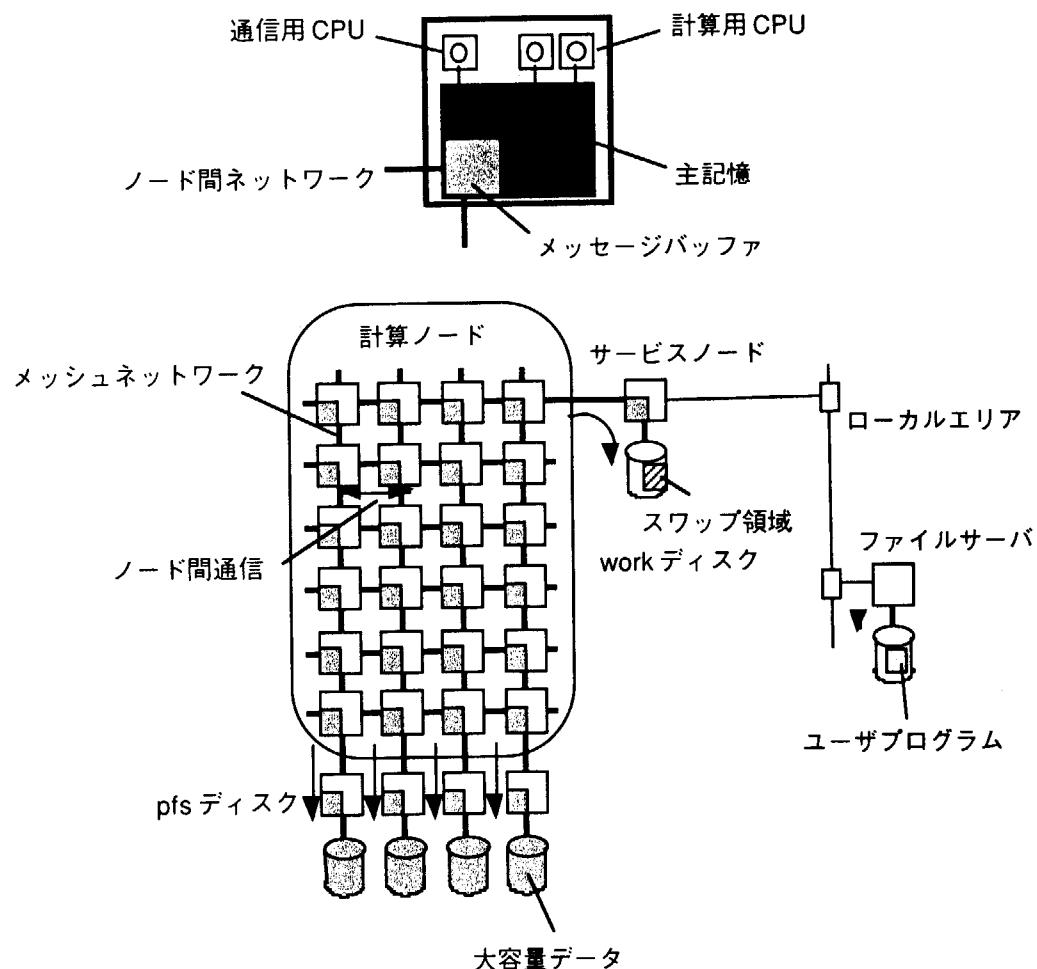


Fig.4.1 Paragon の構成

当然のことだが、入力データや計算された値ですべてのノードで必要となるものは、ノード間で通信して送信しないとデータが共有できない。しかし、この問題は、初歩的だが非常

に重要な問題で、今までシングルプロセッサの計算機しか使ってこなかった者には、相当意識しないとうっかりと忘れてしまうことである。また、データが共有できていて各プロセッサで同じプログラムが走っていたとしても、陽にノード間の同期を取らない限り、同じ命令が同時刻に実行されることは全く保証されない。

一回の通信で転送されるデータの大きさは、メッセージバッファの大きさで制約される。メッセージバッファは、各ノード間ごとの通信用に固定で比例細分化されるので、全バッファの大きさが同じでも、各ノードのバッファ領域はノード数に反比例して減少する仕組みになっているため、大きなデータのノード間の交換や入出力の際に後述のように動作不確実となるプログラムが存在する。

I/O アクセスの問題としては、800 台の計算ノードに対し、ディスクを備えているノードは 20 数台であることから I/O は原理的に低速であることがあげられる。work 領域はサービスノードに接続されているディスクにおかれているため、これに対して集中的にアクセスが行われたとき低速となりシステム障害を起こすことがある。ファイルサーバに接続されたディスクへのアクセスはさらに低速である。

前節で書いたように、通常の ufs とは異なり、Paragon では、大容量データの入出力を高速に行うための pfs が装備され、非常に大きなデータを分割して並列に同時入出力できる。しかし、並列に入出力という概念がシングルプロセッサの計算機にはないため、その使用規約を意識して使わなければならない。

4. Paragon 上での超並列プログラム開発の指針

<一般的注意事項>

1) 超並列計算機は他の汎用計算機やワークステーションに比べてシステム保護機構が弱いため、プログラムを作成するときには、できるだけ “fail to safe” の精神でプログラム設計をする必要がある。

1-1) 第 2 章で説明したとおり、Paragon は、並列計算機ではあるが、個々の計算ノードに UNIX OS が動いているのではなく、個々の計算ノードには、分散 OS (Mach3.0 マイクロカーネル) が搭載され、分散 OS の UNIX に透過的なインターフェイスを介して、サービスノードに搭載されている UNIX OS の機能を利用する構造となっている。超並列計算機は、その最大性能を引き出すために OS レベルでもシステム保護機構が、非常に簡略化されている。このような状況で各 CPU の動作状況を OS に監視させながらプログラムを実行するようにすると、その監視に負荷がかかりすぎて、計算機の実行性能が極端に低下してしまう。したがって、できるだけ並列計算機の性能を引き出すためには、“fail to safe” の精神でプログラム設計をする必要がある。また、現在の並列化コンパイラが汎用計算機のものと比較してまだ未熟であることも大きな原因の一つであるが、当分は、そのような成熟した並列化コンパイラは実現しそうにない。

<プログラム実行>

2) デバッグは、必ず Paragon S5MP で行うこと。また、このことからプログラムを作るときには、ノード数に依存しないソースプログラムを作成しなければならないことになる。もちろん、大きなノード数で生じる問題のデバッグは、この限りではないが、その場合、ラージスケールでのデバッグを開始した旨を情報システム管理課に連絡しておくことを薦める。

2-1) 上記で説明したように、並列計算機では、システム保護機構が非常に簡略化されているため、ユーザのプログラムによってシステム全体をダウンさせる可能性がある。そのため、デバッグは必ず Paragon S5MP で行う。また、並列数に依存するプログラムは、デバッグから本番の計算までの間にプログラムに変更を加える必要が生じるため、デバッグの信頼性を高くするためにもノード数に依存しないソースプログラムを作成しなければならない。さらに、ラージスケールのシミュレーションを行う場合、並列化数の変化によって生じる問題も数多くあるので、Paragon S5MP のデバッグ終了後 Paragon S120MP でもデバッグをしておく必要がある。9) も参照のこと。情報システム管理課の連絡先は、e-mail kanrika@apr.jaeri.go.jp, tel 0720-31-0732 である。連絡は緊急の場合以外、電子メールによって行うこと。

3) 大きなメモリの割り付けは、基本的には動的に行うようにし、もし割り付けられなければ、全ノードが正常にプログラムを終了するようにしておくこと。5)も参照のこと。静的に割り付けると、初期データ（静的割り付けメモリーなど）を全ノードのメモリー上にのせる為、初期データを大量に通信転送しなければならなくなるので、通信やサービスノードに過負荷がかかり、初期ロード時間が非常に長くなる。また、静的に割り付けられたメモリ（ソース、データ、static に割り付けた配列）の大きさは、**size860** コマンドであらかじめ確認しておくこと。各ノードには、128Mbyte のメモリ容量があるが、ユーザの使用するメモリ領域は、通信バッファを含めて 100Mbyte 以下にしておく必要がある。それ以上にするとシステムが不安定になりシステムダウンを引き起こす可能性がある。

3-1) 動的と静的メモリ割り付けのサンプルプログラム

メモリの動的割り付け (dynamic allocate)

サンプルプログラム `memory_dynamic.f` から一部抜粋

```
program memory_dynamic
c
real*4 a
pointer(p,a(200000000))
allocate(a,stat=i)
c
deallocate(a)
stop
end
```

メモリの静的割り付け (dynamic static)

サンプルプログラム `memory_static.f` から一部抜粋

```
program memory_static
c
real*4 a
common /memory/ a
dimension a(20000000)
c
stop
end
```

コンパイル方法

```
% if77 -nx memory_dynamic.f -o memory_dynamic
% if77 -nx memory_static.f -o memory_static
```

3-2) **size860** コマンド

上記コンパイルによって作成したロードモジュールの静的に割り当てた大きさを **size860** コマンドで調べる。一般的には **size** コマンドが知られているが、Paragon 対応として **size860** コマンドを使用した。

text：プログラムリスト自体のサイズ

data：プログラムを実行するにあたっての空きメモリー領域

bss：プログラム内で静的に割り当てたメモリー領域のサイズ

size860 の使用例

```
% size860 -f memory_dynamic
214784(.text) + 37664(.data) + 163488(.bss) = 415936

% size860 -f memory_static
214272(.text) + 37472(.data) + 80163488(.bss) = 80415232
```

ロードモジュール名	text (byte)	data (byte)	bss (byte)
memory_dynamic	214784	37664	415936
memory_static	214272	37472	80163488

3-3) メモリがとれなかった場合の状況の説明

メモリがとれなかった場合、関西研 paragon では標準エラー出力にエラーメッセージが出力される。エラーメッセージとしては以下の 2 つが確認されている。

```
stdio error ; Result too large
Operation not supported on socket
```

これら 2 つのメッセージはいずれも、計算時間が Que-class によって制限された時間を越えるか、もしくは qdel をしたときに標準エラー出力にエラーメッセージとして出力される。つまり、メモリがとれていないプログラムを que に投入すると、関西研の paragon において、S120MP では最大時間まで、S5MP では無制限に時間を費やしてしまうことになる。このため、予め計算時間を把握しておくことが極めて重要となる。

4) 大きなファイルを作るプログラムを実行する場合は、あらかじめ入出力するディスクに十分な空きがあるかを確かめてから実行する。6)でも述べるように、プログラム中でも、ディスクの空きがなくなったとき、I/O処理のエラー値をとり正常に終了できるようにしておく。

4-1) 入出力するディスクの空きの確認方法

入出力するディスクに対して showfs コマンドを発行し、予め容量の確認をしておく。

showfs の使用例 (フラグ k は kbytes 単位を示す。)

```
aprsp% showfs -k /pfs16-2
 Mounted on      kbytes      avail   capacity    sunits     factor
 /pfs16-2       33106560    22060152     26%        65536       16
 sdirs:/home/.sd dirs/vol1h
          /home/.sd dirs/vol2h
          /home/.sd dirs/vol4h
          /home/.sd dirs/vol5h
          /home/.sd dirs/vol7h
          /home/.sd dirs/vol8h
          /home/.sd dirs/vol10h
          /home/.sd dirs/vol11h
          /home/.sd dirs/vol13h
          /home/.sd dirs/vol15h
          /home/.sd dirs/vol16h
          /home/.sd dirs/vol17h
          /home/.sd dirs/vol19h
          /home/.sd dirs/vol20h
          /home/.sd dirs/vol22h
          /home/.sd dirs/vol23h
```

<プログラム設計>

5) プログラム終了時は、全ノードで終了できるように終了前に同期をとる。また、エラーで終了させる場合も、全ノードで終了できるようにエラー処理後、エラー情報を全ノードに送信し、各ノードが同期後にその情報により正常に終了できるようにする。そうしないと、いくつかのノードが終了できないままずっと通信待ちになってしまふ可能性がある。

5-1) 正常にプログラム終了させるプログラムの例 1

ノード 0 番で計算上エラーを起こしてしまい、これ以上の計算が無意味である場合。

サンプルプログラム stop_normal.f

```

program stop_normal
c
include 'fnx.h'
real*8 error_flg,tmp
character*34 check_file
c
error_flg =0.d0
iam = mynode()
c
call gsync()
c error occurred at node 0
if ( iam .eq. 0 ) then
ccc call error_stop(13,dummy,dummy,check_file,error_flg)
ccc same role as call error_stop ccc
error_flg = 1.d0
write(*,*) 'error occurred'
cccccccccccccccccccccccccccccccccccccccc
endif

ccc call error_exit(error_flg)
ccc same role as call error_exit ccc
call gsync()
call gdsum(error_flg,1,tmp)
if( error_flg .gt. 0.d0 ) then
write (*,*) 'error stop'
stop
endif
cccccccccccccccccccccccccccccccccccc
c
call gsync()
stop
end

```

基本的な考えは次の通りである。まず、エラーの有無を `error_flg` という変数に（有り 1、無し 0）で各ノードで蓄えておく。その後で、全ノードの `error_flg` の和をとることにより、全ノードにエラーの有無を伝え、その情報（`error_flg` の和）により、全ノード同期（`gsync`）をとり、終了するようとする。例のプログラムではノード 0 がエラーを発見した（`error_flg = 1.d0`）と想定し、`error_flg` を `gdsum` により和をとることで全ノードで情報を共有し、`error_flg > 0` となるので、全ノードで終了させている。プログラム内の `error_stop`（エラーを発見したとき、そのノードで `error_flg = 1.d0` とする）と `error_exit`（全ノードで `error_flg` の和をとり、その情報からプログラムの終了判定をする）はそれらをツール化したサブルーチンであり、付録 2 の `error.f` に掲載した。

5-2) 正常にプログラム終了させるプログラムの例 2

あるノード（`istop`）に対して過ったエラー処理のあとで、ストップ文を知らずに発行していた場合。

これは、初期の並列化プログラミングにおいて陥り易いプログラミングエラーである。その過った例を下記に示す。

```

error_flg =0.d0
iam = mynode()
c
call gsync()
c
c error occurred at node ix
if ( iam .eq.istop ) then
  error_flg = 1.d0
  write(*,*) 'error occurred'
  stop
endif
c
call gsync()
```

このようなプログラミングをすると全ノードが同期待ち文 `call gsync()` で `ix` ノードの実行を待つため、先にすすめなくなる。このため、`call gsync()` の手前にプログラムの実行中にストップしているノード（`istop`）があれば、全ノードを終了させるツール `call node_check` が必要になる。

```

error_flg =0.d0
iam = mynode()
c
call gsync()
c
c error occurred at node ix
if ( iam .eq.istop ) then
  error_flg = 1.d0
  write(*,*) 'error occurred'
  stop
endif
c
call node_check(error_flg)
c
call gsync()

```

サンプルプログラム内の **node_check**. はノード間の非同期通信命令を使い、ストップしているノードがあれば、他のノードの全てを終了させるサブルーチンであり、付録2の node_check.f に掲載した。

6) ファイル処理には、エラー処理をいれる。ファイル処理を行ったノードでエラーが起こり、そのノードだけが終了してしまうときがある。この場合、ジョブの実行状態を端末から確かめると正常に実行を継続をしているように見えるが、実際には、1 ノードが終了し、他のノードは通信待ちになっていてプロセスは全く進んでいない。

6-1) I/O 処理のエラー値をとり正常に終了させるプログラム。

サンプルプログラム file_access.f

```

program file_access
c
include 'fnx.h'
integer iam
real*8 error_flg
character*34 check_file
c
call gsync()
iam = mynode()

```

```

c
check_file='test_file'
call file_exist_check(check_file,error_flg)
if (iam .eq. 0) then
  call sopen(10,check_file,error_flg)
  write (10,*) 'open test'
endif
close(10,status='keep')
call error_exit(error_flg)
c
stop
end

```

このプログラムはファイルアクセスのうち **open** についてサブルーチン化した (**call sopen**) を使用したものである。**sopen** は、オープン時にエラーがあると **error_flg=1** を代入するサブルーチンであり、付録2の **sopen.f** に掲載した。I/O処理のエラー値をとり正常に終了させることは関西研 Paragon ではファイルのアクセス (**open**, **read**, **write**, **inquire** 等) にはノード数制限があるために、並列化プログラミングにおいて必要不可欠なことである。これについては詳細を 9)に記述する。ここで、サブルーチン **error_exit** は、前述したように全ノードで **error_flg** の和をとり、それが 0 より大きければ全ノード終了させる働きをしている。

7) プログラム中で UNIX コマンドをシステム文を用いて呼び出すことは極力さけること。システム文で UNIX コマンドを発行すると、これはサービスノードにおいて子プロセストシテ実行されるようである。UNIX コマンドは元のプログラムと非同期で実行されるので、UNIX コマンドで行った **mkdir** 等の処理が完了強いたことを **inquire** 等の別のコマンドで確認することが必要である。また、同じ UNIX コマンドを多ノードで同時に発行することはシステムの負荷を非常に高めるので避けなければならない。例えば、システム文を用いて **mkdir** の発行の後、そのディレクトリ下にファイルを作ろうとする場合、ディレクトリがまだできていないのにファイルを作ろうとして、エラーとなる可能性がある。

7-1) **call system** (C 言語の関数) と **i=system** (Fortran のシステム関数) の例とそれと等価なシェルの例

サンプルプログラム **mkdir_system1.f**

```

program mkdir_system1
c
include 'fnx.h'
integer iam
c
iam = mynode()
c
if (iam .eq. 0 ) then
  call system('mkdir test_dir')
  open (10, file= 'test_dir/data')
  write (10,*) 'hallo'
  close (10,status='keep')
endif
c
stop
end

```

サンプルプログラム mkdir_system2.f

```

program mkdir_system2
c
include 'fnx.h'
integer iam,i_system
c
iam = mynode()
c
if (iam .eq. 0 ) then
  i_system=system('mkdir test_dir')
  open (10, file= 'test_dir/data')
  write (10,*) 'hallo'
  close (10,status='keep')
endif
c
stop
end

```

上記サンプルプログラムではファイルのオープンの直前に system 文もしくは system 関

数を使い親ディレクトリを作成している。ファイルをオープンする前に親ディレクトリを作成すること自体は正しいことであるが、シェルの中で作成しておけば、より信頼性の向上と高速化につながる。下記にそのシェルスクリプトを示す。

```
#!/bin/csh
# qsub and mkdir

qsub -r test          ¥
-o ./out.log          ¥
-e ./err.log          ¥
-q sppss4-i << EOF

if (! -d test_dir ) then
    mkdir test_dir
endif

ロードモジュール名 -plk -sz 4 -mbf 10240000

exit
EOF
```

<ファイル処理関係>

8) ファイル処理する直前に必ずその親ディレクトリの存在を 1 ノードに制限した inquire 文を使って確かめ、ファイルが存在しない場合は、正常終了させるようにプログラミングする。これを怠ると、プロセスを握ったままの状態になり、qdel してもそのプロセスは解放されず、いくつかのノードが新たなプロセスを受け付けなくなり、Paragon を再起動しなければならなくなる。

8-1) inquire 文の例。そのときの qdel の方法。

先の紹介したサンプルプログラム中にある

```
check_file='test_file'
call file_exist_check(check_file)
サブルーチン file_exist_check が該当する。
サンプルプログラム file_exist.f
```

```

subroutine file_exist_check(check_file)
c
  include 'fnx.h'
c
  integer iam
  character*34 check_file
  logical*4 ex
c
  iam = mynode()
  if (iam .eq. 0) then
    inquire(file=check_file,exist=ex)
  endif

  if ( .not. ex ) then
    call error_stop(18,idummy,idummy,check_file.error_flg)
  end if
c
  call error_exit(error_flg)
c
  return
end

```

このサブルーチン `file_exist_check` では変数名 `check_file` に割り当てられたファイルもしくはディレクトリの存在チェックを `inquire` 文で 1 ノードに制限して行っている。ファイルもしくはディレクトリの存在・不在の情報（`ex` の値）を使って、プログラムを終了させている。

9) 1 ノードがファイルを開くことができる数（論理機番を割り付けられる数）は 61 個である。また、1 つのディスクに対して同時に `open`、`read` (`write`) 入出力のできるノード数は、32 未満とする。32 ノード以上で同時に入出力するとシステムダウンを引き起こす可能性がある。

9-1) 論理機番を割り付けられる数とノード数の制限

論理機番を割り付けられる数は実際には 64 個であるが、そのうち標準エラー出力 0 番、標準入力 5 番、標準出力 6 番とすでに割り付けられているため、61 個である。

ノード数の制限には、2 通りある。1 つはプログラムの中で分業するような作業がある場合に計算を行うノードを制限する場合と、ファイルにアクセスするためにノードを制限する場合がある。そのどちらにおいてもエラーが起きた処理を全ノードで正しくする必要

がある。

10) 32 ノード以上で同時に入出力するためには、**gopen** という命令を使い、ファイルを開き、**cread (cwrite)**、**iread (iwrite)** を使って入出力しなければならない。ただし、**cread (cwrite)** は、同期型並列入出力であり、詳細なエラー情報がとれない為、できるだけ非同期型並列入出力の **iwrite (iread)** を使い、入出力の確実な終了確認をとり、エラー時も正常にプログラムを終了させなければならない。これらの並列入出力は、オプション **M_RECORD** でストライピング（イメージは1つだが物理的には複数のディスクに分けられて格納されている）された pfs ファイルとして pfs ディスクにのみ並列入出力ができる。同様に、**M_SYNC** でストライピングしないファイルとして ufs ディスクに入出力ができる。

10-1) gopen M_record , cread +エラー処理の例。

サンプルプログラム `cread_M_RECORD.f`

```

program cread_M_record
c
  include 'fnx.h'
  integer iam,nodes
  real*8 error_flg
  character*34 check_file
c
  iam = mynode()
  nodes= numnodes()
  call gsync()
  check_file='/pfs16-2/j6446/'
  call file_exist_check(check_file,error_flg)
c
  call gopen(10,check_file//'test',M_RECORD)
  if (iomode(10).ne.M_RECORD) then
    close(10)
    call error_stop(17,dummy,dummy,check_file,error_flg)
  end if
  call error_exit(error_flg)
c
  gsync()
c

```

```

irw = (nodes+15)/16
do ir = 1,irw
  if ( imod(iam,irw).eq.ir-1 ) then
    call cread(10, nodes, 4)
  endif
enddo
close(10)
c
stop
end

```

上記プログラムの **gopen** のモード **M_RECORD** は、パラゴンで最も早い入出力方法である。内部的な処理として、個別のファイルポインタをもち、全ノードともノンブロッキングで同じオペレーションを同じ順番でおこなう。このプログラムの流れは次の通りである。まず、サブルーチン **file_exist_check** で入出力するファイルが存在するかどうか確認し、**gopen** でファイルをオープンする。次に、**iomode** によって正常に **gopen** できたかどうかを確かめ、次に **gsync** で全ノード同期をとる。最後に、より高速に並列入出力をを行うためにストライプ数と同数のノードずつ **cread** を実行する。このサンプルプログラムは、**cread** を通るノードが常に 16 ノードに制限されていることから、16 ストライピングされた pfs ファイルに対して有効的な手法をとっているといえる。

10-2) **gopen M_sync , cread +エラー処理の例。**

サンプルプログラム **cread_M_SYNC.f**

```

program cread_M_SYNC
c
include 'fnx.h'
integer iam,nodes
real*8 error_flg
character*34 check_file
c
iam = mynode()
nodes = numnodes()
call gsync()
check_file='/pfs16-2/j6446/'
call file_exist_check(check_file, error_flg)
c

```

```

call gopen(10,check_file //'test',M_SYNC)
if (iomode(10).ne.M_SYNC) then
  close(10)
  call error_stop(17,dummy,dummy,check_file,error_flg)
end if
call error_exit(error_flg)

c
call cread(10, nodes, 4)
close(10)

c
stop
end

```

サンプルプログラムの流れは次の通りである。まず、サブルーチン `file_exist_check` で入出力するファイルが存在するかどうか確認し、`gopen` でファイルをオープンする。次に、`iomode` によって正常に `gopen` できたかどうかを確かめ、`cread` を実行する。上記プログラムの `gopen` のモード `M_SYNC` は、内部的な処理として、全てのノードが同じファイルポインタを持つため、並列入出力向きではないといえる。利用方法としてストライピングしないファイルとして ufs ディスクに入出力することが考えられる。ファイルにアクセスに行く順番がノード番号順であるために、容量の大きいファイルの入出力すると時間がかかることになり、そのあいだは、ノードごとの逐次処理になるため並列化効率の劣化を招くことになる。しかし、これらのプログラムの `cread` (`cwrite`) は同期型並列入出力であり、詳細なエラー情報がとれない為、入出力の確実な終了確認ができない。そのため、入出力の確実な終了確認をとり、エラー時も正常にプログラムを終了できる、非同期型並列入出力の `iwrite` (`iread`) を使うほうがよい。

10-3) `gopen M_record , iwrite +エラー処理の例。`

サンプルプログラム `iwrite_M_RECORD.f`

```

program iwrite_M_record
c
include 'fnx.h'
integer iam,nodes,iwrite_id
real*8 error_flg,start_time,total_time
character*34 check_file
c
iam = mynode()

```

```

nodes = numnodes ()
call gsync()
check_file='/pfs16-2/j6446/'
call file_exist_check(check_file,error_flg)

c
call gopen(10,check_file//'test',M_RECORD)
if (iomode(10).ne.M_RECORD) then
    close(10)
    call error_stop(17,dummy,dummy,check_file)
end if
call error_exit(error_flg)

c
irw = (nodes +15)/16
do ir = 1,irw
    if ( imod(iam,irw).eq.ir-1 ) then
        iwrite_id =iwrite(10, nodes, 4)
        total_time = 0.d0
        start_time = dclock()
        do while(iodone(iwrite_id) .eq. 0)
            total_time = dclock() - start_time
            if ( total_time .gt. 1.d1 ) then
                call error_stop(14,dble(iam),dummy,error_flg)
            endif
        enddo
    endif
    call gsync()
    call error_exit(error_flg)
enddo
close(10)

c
stop
end

```

非同期型並列入出力の **iwrite** を使用した例である。基本的には、前の同期入出力と同じだが、**iodone** で完了ステータスが取れることを利用して、**iodone** 完了ステータスの合否を判定条件にして無限ループ構造をつくり、長い時間ループを抜けられない時「エラーが起こった」と判断することにしている。

```

if ( total_time .gt. 1.d1 ) then
    call error_stop(14,dble(iam),dummy,check_fileerror_flg)
endif

```

の1行目で指定するべき時間（下線部）は write される変数の大きさによって変わるので注意が必要である。

11) pfs ファイルは、通常の ufs ファイルではないため非フォーマットでしかデータを保存できなく、また、常にそのデータは作成したときと同じノード数の並列プログラムでしか利用できない。さらに、2Gbyte (32bit アドレッシングの上限値) 以上のファイルを作成することが可能であるが、2Gbyte 以上のファイルに対しては、ほとんどの (**ls**、**mv**、**cat** などは有効であるが) 従来の UNIX コマンドを受け付けなくなる。

11-1) 2Gbyte 以上のファイルが判定できなかったシェルの例。

```

#!/bin/csh
if ( -f restart_file ) then
    echo "ok"
    qsub restart.bat
else
    echo "ng"
endif
exit

```

シェルスクリプトを組んでジョブのリストアート等を行う上で注意するべき点である。上記のシェルスクリプトの場合 `restart_file` が存在すればリストアート用の `qsub` を行うことになる。ここで問題なのは、`restart_file` のサイズである。サイズが 2Gbyte 以上である場合には、判定 `if (-f restart_file)` が出来ないためにリストアートが出来なくなってしまう。

<通信関係>

12) ノード間の情報の送受信は、**csend** (**crecv**)、**isend** (**irecv**) で行うが、双方向送受信時では **csend** と **crecv** の組み合わせは、システムデッドロックを引き起こす可能性があるので使用は避けなければならない。また、**isend** と **irecv** 組み合わせで

は、送受信が正常に完了したかの返値を message 関数に入れて必ず確認する。**isend** と **crecv** を組み合わせることにより、システムデッドロックを引き起こさないだけでなく、終了確認も省略する簡易的な使い方ができる。

12-1) isend crecv +エラー処理の例。

サンプルプログラム `tarns_check.f`

```

subroutine trans_check
1  (trans_type,send_buf,send_count,send_node,send_ptype,
2           recv_buf,recv_count,error_flg)

c
include 'fnx.h'
character*34 check_file

c
integer trans_type,send_count,send_node,send_ptype,
1      recv_count,trans_id,trans_sum_id,iam, nodes,counter,
2      isend_id,irecv_id
byte   send_buf(send_count),recv_buf(recv_count)
real*8 wait_t_time,start_t_time,error_flg

c
iam = mynode()
nodes = numnodes ()
trans_id = 0
counter = 0
wait_t_time = 0.0
c  if (iam.eq.0) then
c    write (*,97) trans_type
c97  format('subroutine trans_check',i5,' start')
c  endif
call gsync()
start_t_time = dclock()

isend_id = isend
1           (trans_type,send_buf,send_count,send_node,send_ptype)
c  call csend (trans_type,send_buf,send_count,send_node,send_ptype)
c
call crecv (trans_type,recv_buf,recv_count)
do while( trans_id .eq. 0 )

```

```

        counter = counter + 1
        trans_id = msgdone(isend_id)
        wait_t_time = dclock() - start_t_time
        if( wait_t_time .gt. 60d0 ) then
            trans_id = 1
        endif
    end do
c   write (*,98) iam,counter,wait_t_time
c   call msgwait( isend_id )
c
c       call gsync()
c       call gisum( trans_id , 1 , tmp )
c98  format ('iam=',i3,' counter=',i5,' wait_t_time=',1p,e15.8)
c       call gsync()
        if (trans_id .ne. nodes) then
            call error_stop(20,idummy,idummy,check_file,error_flg)
        endif
        call error_exit(error_flg)
c       if (iam.eq.0) then
c           write (*,99) trans_type,trans_id
c99  format('subroutine trans_check',i5,' end trans_id=',i3)
c       endif
c
        return
    end

```

上記サブルーチンは、双方向送受信時に **csend** と **crecv** の組み合わせを用いるとシステムデッドロックを引き起こす可能性がある（付録3参照）ために **isend** と **crecv** の組み合わせにして、システムデッドロックを回避しているプログラムである。

13) 通信のバッファサイズは、実行時のオプション **-mbf ???** （???はバッファサイズ [byte]）で行う。一回の通信で転送されるデータの大きさは、メッセージバッファの大きさで制約される。メッセージバッファは、各ノード間ごとの通信用に固定で比例細分化（付録3参照）されるので、全バッファの大きさが同じでも、各ノードのバッファ領域はノード数に反比例して減少する仕組みになっているため、注意が必要である。

13-1) **csend crecv** のデッドロックプログラムと PE 数依存性。
前項 12) で用いたサブルーチン **trans_check** において **isend** を **csend** と置き換える。

通信量をふやし、**-mbf ???**を減らすとシステムデッドロックを引き起こす。詳細については付録3 メッセージパッシングの手法に記述する。

<コンパイル・実行時のオプション>

14) ロードモジュールのコンパイルの最適化オプションは、Intel に相談後に既定値から変更すること。**-nx** オプションは必ずつけること。**-nx** オプションとは、プログラムを計算ノード上で実行させるためのもので、これがなければサービスノード上で実行され、I/O のレスポンスの悪化などを招く。特に、デバック中は**-Mbounds** (配列の領域外参照チェック) をつけておく。

15) 実行時のオプションには、**-p1k** オプションは必ずつけること。**-p1k** オプションとは、メモリ上に載らなかった情報をいったんディスクに待避させていわゆる仮想記憶を使用にするオプションである。このオプションを適切に作動させるためには、メモリを割り付けた直後、常に全ノード通信を行いメモリをロックする必要がある。これをしなければ、Paragon では全てのノードがディスクにつながっているわけではないので仮想記憶時に非常に大きな通信負荷が生じ、実質的に計算速度が激減する。

16) Nx ライブラリーを使用するサブルーチンやプログラムでは、**fnx.h** をインクルードしておくこと。

5. 結び

Paragon の動作不確実の原因の一つは、ディスク I/O が通信負荷を増加させることにある。各ノードがディスクを持ち、スワップ領域をローカルディスク上にとれば、通信負荷を増やす、仮想記憶を活用することができてプログラムの柔軟性が増えるだろう。また、メッセージバッファ容量に関して、初期に指定した大きさをノード数で均等に分ける方式はシステムスケーラビリティを損なっており、より柔軟にメッセージバッファを活用できることが望まれる。さらに、FORTRAN の通信及び入出力システム関数などの関数やサブルーティンがより適切なエラーを返すことが、プログラムの信頼性を高めるために望まれる。

謝辞

本研究は、日本原子力研究所関西研究所光量子科学センターで行われたものであり、すばらしい研究の場を提供していただくとともに、研究に対する深い御理解と御支援をいただいた大野英雄所長に謝意を表します。協力体制を整備し研究を推進していただきました有澤孝センター長に深く感謝いたします。また、光量子シミュレーション研究グループにおいて、この研究を強く御支援していただきました田島俊樹グループリーダー、井原均サブグループリーダー、那珂研究所炉心プラズマ研究部プラズマ理論研究室長岸本泰明氏に深厚なる謝意を表します。最後に、関西研駐在の情報システム管理課及びシステムエンジニア、Intel の皆様には、いつも計算機運用支援及び情報提供などでお世話になりました。改めて、ここで厚くお礼申し上げます。

参考文献

- 1) 関西研究所 並列計算機 (Paragon) 利用手引き
日本原子力研究所計算科学技術推進センター
- 2) PARAGON TM Fortran Compiler User's Guide
- 3) PARAGON TM Fortran System Calls Reference Manual
- 4) 佐々木明 “MPI を用いた 2 次元流体シミュレーション”，
JAERI-Data/Code 97-048
- 5) K. Tani, "Advanced Phton Simulation",
Journal of Plasma and Fusion Reserch 72, 935 (1996)
(光量子シミュレーション、プラズマ・核融合学会誌)

付録1 サンプルプログラム メインルーチン

cread_M_RECORD.f	: オプション M_RECORD の同期並列入力
cread_M_SYNC.f	: オプション M_SYNC の同期並列入力
cwrite_M_RECORD.f	: オプション M_RECORD の同期並列出力
cwrite_M_SYNC.f	: オプション M_SYNC の同期並列出力
file_access.f	: 入出力処理のエラー値をとり異常があった場合、正常にプログラムを終了
iwrite_M_RECORD.f	: オプション M_RECORD の非同期並列出力
memory_dynamic.f	: 動的なメモリー割り付け
memory_static.f	: 静的なメモリー割り付け
mkdir_system1.f	: システムコールによるディレクトリーの作成
mkdir_system2.f	: システム文によるディレクトリーの作成
stop_abnormal.f	: あるノードが予期せず終了していたときの全ノード終了
stop_normal.f	: あるノードでエラーを検知して全ノード終了
test_trans.f	: 通信が正常に終了したかどうかを確認

```

c=====
c      program cread_M_record
c          error.f is needed
c          file_exist.f is needed
c=====
c      program cread_M_record
c      include 'fnx.h'
c      integer iam, nodes
c      real*8 error_flg
c      character*34 check_file
c
c      iam = mynode()
c      nodes = numnodes ()
c      check_file='pfs16-2/j6446/'
c      call gsync()

c      checking the existence of a file
c      call file_exist_check(check_file,error_flg)
c
c      confirmation of fopen at option= M_RECORD
c      call fopen(10,check_file//'test',M_RECORD)
c      if (iomode(10).ne.M_RECORD) then
c          close(10)
c          call error_stop
c          &      (17,dummy,dummy,check_file,error_flg)
c      end if
c      call error_exit(error_flg)

c      16 splitting cread
c      irw = (nodes +15)/16
c      do  ir = 1,irw
c          if ( imod(iam,irw).eq.ir-1 ) then
c              call cread(10, nodes, 4)
c          endif
c      enddo
c      close(10)
c      stop
c      end

```

```

c=====
c      program cread_M_SYNC
c              error.f is needed
c              file_exist.f is needed
c=====
c
c      program cread_M_SYNC
c          include 'fnx.h'
c          integer iam, nodes
c          real*8 error_flg
c          character*34 check_file
c
c          iam = mynode()
c          nodes = numnodes ()
c          check_file='pfs16-2/j6446/'
c          call gsync()
c
c          checking the existence of a file
c          call file_exist_check(check_file)
c
c          confirmation of fopen at option= M_SYNC
c          call fopen(10,check_file//'test',M_SYNC)
c          if (iomode(10).ne.M_SYNC) then
c              close(10)
c              call error_stop
c              &      (17,dummy,dummy,check_file, error_flg)
c          end if
c          call error_exit(error_flg)
c
c          call cread(10, nodes, 4)
c          close(10)
c
c          stop
c      end

```

```

c=====
c      program cwrite_M_record
c          error.f is needed
c          file_exist.f is needed
c=====
c
c      program cwrite_M_record
c      include 'fnx.h'
c      integer iam, nodes
c      real*8 error_flg
c      character*34 check_file
c
c      iam = mynode()
c      nodes = numnodes ()
c      check_file='/pfs16-2/j6446/'
c      call gsync()
c
c      checking the existence of a file
c      call file_exist_check(check_file)
c
c      confirmation of fopen at option= M_RECORD
c      call fopen(10,check_file//'test',M_RECORD)
c      if (iomode(10).ne.M_RECORD) then
c          close(10)
c          call error_stop
c          &      (17,dummy,dummy,check_file, error_flg)
c      end if
c      call error_exit(error_flg)
c
c      16 splitting cwrite
c      irw = (nodes +15)/16
c      do ir = 1,irw
c          if ( imod(iam,irw) .eq. ir-1 ) then
c              call cwrite(10, nodes, 4)
c          endif
c      enddo
c      close(10)
c      stop
c  end

```

```

c=====
c      program cwrite_M_SYNC
c              error.f is needed
c              file_exist.f is needed
c=====
c
c      program cwrite_M_SYNC
c          include 'fnx.h'
c          integer iam, nodes
c          real*8 error_flg
c          common /error$/ error_flg
c          character*34 check_file
c
c          iam = mynode()
c          nodes = numnodes ()
c          check_file='/pfs16-2/j6446/'
c          call gsync()
c
c          checking the existence of a file
c          call file_exist_check(check_file)
c
c          confirmation of fopen at option= M_SYNC
c          call fopen(10,check_file//'test',M_SYNC)
c          if (iomode(10).ne.M_SYNC) then
c              close(10)
c              call error_stop
c              &      (17,dummy,dummy,check_file, error_flg)
c          end if
c          call error_exit(error_flg)
c
c          call cwrite(10, nodes, 4)
c          close(10)
c
c          stop
c      end

```

```
c=====
c      program file_access
c              error.f is needed
c              file_exist.f is needed
c              sopen.f is needed
c=====
c
c      program file_access
c      include 'fnx.h'
c      integer iam
c      real*8 error_flg
c      character*34 check_file
c
c      iam = mynode()
c
c      check_file='test_file'
c      call gsync()
c
c      checking the existence of a file
c      call file_exist_check(check_file)
c      if (iam .eq. 0) then
c          call sopen(10, check_file,error_flg )
c          write (10,*) 'open test'
c      endif
c      close(10,status='keep')
c      call error_exit(error_flg)
c
c      stop
c      end
```

```

c=====
c      program iwrite_M_record
c          error.f is needed
c          file_exist.f is needed
c          wait.f is needed
c=====
c      program iwrite_M_record
c      include 'fnx.h'
c      integer iam, nodes,iwrite_id
c      real*8 error_flg,start_time,total_time
c      character*34 check_file

c
c      iam = mynode()
c      nodes = numnodes ()
c      check_file='/pfs16-2/j6446/'
c      call gsync()

c      checking the existence of a file
c      call file_exist_check(check_file)

c
c      confirmation of fopen at option= M_RECORD
c      call fopen(10,check_file//'test',M_RECORD)
c      if (iomode(10).ne.M_RECORD) then
c          close(10)
c          call error_stop
c          &      (17,dummy,dummy,check_file, error_flg)
c      end if
c      call error_exit(error_flg)

c
c      confirmation of 16 spliting iwrite
c      irw = (nodes +15)/16
c      do  ir = 1,irw
c          if ( imod(iam,irw).eq.ir-1 )  then
c              iwrite_id =iwrite(10, nodes, 4)
c              total_time = 0.d0
c              start_time = dclock()

```

```
do while(iodone(iwrite_id) .eq. 0)
    total_time = dclock() - start_time
    if ( total_time .gt. 1.d1 ) then
        call error_stop
&        (14,dble(iam),dummy,check_file,error_flg)
        endif
    enddo
    endif
    call gsync()
    call error_exit(error_flg)
enddo
close(10)

c
stop
end
```

```
c=====
c      program memory_dynamic
c=====

program memory_dynamic
include 'fnx.h'
real*4 a

pointer(p,a(200000000))
allocate(a,stat=i)

c
iam = mynode()
c
a(1) = 1.0d0
a(123456789) = 2.0d0
c
write (*,*) iam,i,a(1),a(123456789)
100 format('iam=',i2,' stat=',i1,2f15.8)
c
call fopen(51,'/work01/j6446/data01',M_SYNC)
call cwrite (51,real(iam),4)
close (51,status='keep')
c
deallocate(a)
stop
end
```

```
C=====
c      program memory_static
C=====

      program memory_static
      real*4 a
      include 'fnx.h'
      common /memory/ a
      dimension a(20000000)

c
      iam = mynode()
c
      a(1) = 1.0do
      a(123456789) = 2.0do
c
      write (*,*) iam,i,a(1),a(123456789)
100   format('iam=',i2,' stat=',i1,2f15.8)
c
      call gopen(51,'/work01/j6446/data01',M_SYNC)
      call cwrite (51,real(iam),4)
      close (51,status='keep')
c
      call gsync()

      stop
      end
```

```
c=====
c      program mkdir_system1
c=====

program mkdir_system1
include 'fnx.h'
integer iam

c
iam = mynode()

c
if (iam .eq. 0 ) then
  call system('mkdir test_dir')
  open (10, file= 'test_dir/data')
  write (10,*) 'hallo'
  close (10,status='keep')
endif

c
stop
end
```

```
c=====
c      program mkdir_system2
c=====

program mkdir_system2
include 'fnx.h'
integer iam,i_system

c
iam = mynode()

c
if (iam .eq. 0 ) then
  i_system=system('mkdir test_dir')
  open (10, file= 'test_dir/data')
  write (10,*) 'hallo'
  close (10,status='keep')
endif

c
stop
end
```

```
c=====
c      program stop_abnormal
c              node_check.f is needed
c              error.f is needed
c=====
c      program stop_abnormal
c          include 'fnx.h'
c          real*8 error_flg
c
c      iam = mynode()
c
c      node 1 stop
c      if ( iam .eq. 1 ) then
c          write (*,*) 'iam=',iam,' is stop'
c          stop
c      endif
c
c      checking of stopping nodes
c      call node_check(error_flg)
c
c      stop
c      end
```

```

c=====
c      program stop_normal
c                      error.f is needed
c=====

      program stop_normal
      include 'fnx.h'
      real*8 error_flg,tmp
      character*34 check_file

c
      error_flg =0.d0
      iam = mynode()

c
      call gsync()

c
c      error occurred at node 0
      if ( iam .eq. 0 ) then
          error_flg = 1.d0
          write(*,*) 'error occurred'
cccc      call error_stop
cccc &           (13,dummy,dummy, check_file, error_flg)
      endif

c
c      transference of an error information and all node stop
      call gsync()
      call gdsum(error_flg,1,tmp)
      if( error_flg .gt. 0.d0 ) then
          write (*,*) 'error stop'
          stop
      endif

c
cccc      call error_exit(error_flg)
c
      call gsync()
      stop
end

```

```

c=====
c      program test_trans
c              trans_check.f is needed
c              error.f is needed
c=====
c
c      program test_trans
c      include 'fnx.h'
c
c      integer is,iam, nodes,nid,isend_id,n_time
c      real*8 a,b,error_flg,tmp
c      character*34 check_file
c      dimension a(150000),b(150000)
c
c      do i= 1,150000
c          a(i) = 0.0
c          b(i) = 0.0
c      enddo
c
c      a(1)=1.0
c
c      do i = 2, 150000
c          a(i) = i + a(i-1)
c      enddo
c
c      call gsync()
c      iam = mynode()
c      nodes = numnodes ()
c
c      do i= 1,150000
c          a(i) = iam * a(i)
c      enddo
c
c      c initial data at node 0
c      if (iam.eq.0) then
c          write (6,10) (i,a(i),i=1,10)
c          write (6,10) (i,a(i),i=14991,15000)
c      endif

```

```

c
call gsync()

c
c..... set the "network ID"
if ( iam .eq. nodes -1 ) then
    nid = 0
else
    nid = iam + 1
endif

c
c..... send & recieve the buffer data
c   (nodes - 1) times transference to next node
n_time = nodes - 1
do i=1,n_time
    if (iam.eq.0) then
        write (*,*) 'loop = ',i
    endif
    call trans_check
1    (777,a,8*150000,nid,0,
2          b,8*150000,
3          check_file,error_flg )
    enddo

c
c   data after transference at node (nodes - 1)
if (iam.eq.nodes) then
    write (6,20) (i,is,isend_id,b(i),i=1,10)
    write (6,10) (i,b(i),i=14991,15000)
endif
10   format (i2,1x,1p,e15.8)
20   format (3i2,1x,1p,e15.8)

c
call gsync()

stop
end

```

付録2 サンプルプログラム サブルーチン

error.f	: エラーフラグへの値代入とエラーメッセージ出力
error_exit.f	: エラーフラグの全ノード和と全ノード終了
file_exist.f	: ファイルの存在確認
node_check.f	: ノードの動作-非動作の確認
sopen.f	: エラー処理付きファイルオープン
trans_check.f	: 通信完了確認
wait.f	: 時間待ち

```

c=====
c      subroutine error_stop
c          checking errors
c=====
subroutine error_stop(no_err,prm1,prm2,err_file, error_flg)
c
integer no_err
real*8 prm1,prm2
character*34 err_file
c
error_flg = 1.d0
if( no_err .eq. 13 )  then
  write(*,*) 
  '### error 13 :, ' error'
end if
if( no_err .eq. 14 )  then
  write(*,*) 
  '### error 14 : iam=',prm1,' iwrite error'
end if
if( no_err .eq. 16 )  then
  write(*,*) 
  '### error 16 : ',err_file,' cannot open'
end if
if( no_err .eq. 17 )  then
  write(*,*) 
  '### error 17 : ',err_file,' is iomode error'
end if
if( no_err .eq. 18 )  then
  write(*,*) 
  '### error 18 : ',err_file,' is none'
end if
if( no_err .eq. 20 )  then
  write(*,*) 
  '### error 20 : ', ' isend-crecv error'
end if

return
end

```

```
c=====
c      subroutine error_exit
c      program stop if error occurred
c=====
c      subroutine error_exit(error_flg)
c
c      include 'fnx.h'
c
c      call gdsum(error_flg,1,tmp)
c      if( error_flg .gt. 0.d0 ) then
c          if (iam .eq. 0) then
c              open (99,file='err.out', status='unknown')
c              write (99,*) 'stopped by error_exit'
c              close (99)
c          endif
c          call gsync()
c          stop
c      endif
c
c      return
c  end
```

```
c=====
c      subroutine file_exist_check
c              error.f is needed
c=====
c      subroutine file_exist_check(check_file,error_flg)
c
c      include 'fnx.h'
c
c      integer iam, nodes
c      character*34 check_file
c      logical*4 ex
c
c      iam = mynode()
c      if (iam .eq. 0) then
c          inquire(file=check_file,exist=ex)
c      endif
c
c      if ( .not. ex ) then
c          call error_stop
c          &      (18,idummy,idummy,check_file,error_flg)
c      end if
c
c      call error_exit(error_flg)
c
c      return
c      end
```

```

c=====
c      subroutine node_check
c          wait.f is needed
c=====
c      subroutine node_check
include 'fnx.h'
integer iam,isend_id,irecv_id,check_node
1      , trans_send_id(0:511),trans_recv_id(0:511)
real*8 total_time,start_time
1      , wait_t_time , start_t_time

c
iam = mynode()
node = numnodes ()

c
call wait_time(1.d0)

c
total_time = 0.d0
start_time = dclock()

c
do check_node = 0, nodes-1
    wait_t_time = 0.d0
    start_t_time = dclock()
    trans_send_id(check_node) = 0
    trans_recv_id(check_node) = 0

c
if (iam .eq. check_node) then
    isend_id = isend(check_node,total_time,8,-1,0)
    do while(trans_send_id(check_node) .eq. 0)
        trans_send_id(check_node) = msgdone(isend_id)
        wait_t_time = dclock() - start_t_time
        if ( wait_t_time .gt. 5.d-1 ) then
            total_time = dclock() - start_time
            write (*,*) 'iam=',iam
            1      , ' is stoped by isend error',total_time
            stop
        endif
    enddo

```

```
else
    irecv_id = irecv(check_node,total_time,8)
    do while(trans_recv_id(check_node) .eq. 0)
        trans_recv_id(check_node) = msgdone(irecv_id)
        wait_t_time = dclock() - start_t_time
        if ( wait_t_time .gt. 5.d-1 ) then
            total_time = dclock() - start_time
            write (*,*) 'iam=',iam
            1      , ' is stoped by irecv error',total_time
            stop
        endif
    enddo
    endif
    call wait_time(5.d-1)
enddo

c
total_time = dclock() - start_time
write (*,*) 'iam=',iam,' is not stoped by node_check'
1      ,total_time

c
return
end
```

```
c=====
c      subroutine sopen
c              error.f is needed
c=====
subroutine sopen(logical_unit,file_name,error_flg)

character*34 file_name

c
open ( logical_unit, file=file_name, err=100 )
c
return
100 call error_stop
&      (16,idummy,idummy,file_name,error_flg)
c
return
end
```

```

c=====
c      subroutine trans_check
c          error.f is needed
c=====

      subroutine trans_check
1  (trans_type,send_buf,send_count,send_node,send_ptype,
2           recv_buf,recv_count,
3   check_file,error_flg)

c
      include 'fnx.h'
      character*34 check_file
      real*8 error_flg

c
      integer trans_type,send_count,send_node,send_ptype,
1       recv_count,trans_id,trans_sum_id,iam, node,counter,
2       isend_id,irecv_id
      byte   send_buf(send_count),recv_buf(recv_count)
      real*8  wait_t_time,start_t_time

c
      iam = mynode()
      nodes = numnodes ()
      trans_id = 0
      counter = 0
      wait_t_time = 0.0
c      if (iam.eq.0) then
c          write (*,97) trans_type
c97      format('subroutine trans_check',i5,' start')
c      endif
      call gsync()
      start_t_time = dclock()

      isend_id = isend
1           (trans_type,send_buf,send_count,send_node,send_ptype)
c      call csend (trans_type,send_buf,send_count,send_node,send_ptype)
c
      call crecv (trans_type,recv_buf,recv_count)
      do while( trans_id .eq. 0 )

```

```
counter = counter + 1
trans_id = msgdone(isend_id)
wait_t_time = dclock() - start_t_time
if( wait_t_time .gt. 60d0 ) then
    trans_id = 1
endif
end do
c   write (*,98) iam,counter,wait_t_time
c   call msgwait( isend_id )
c
c   call gsync()
call gisum( trans_id , 1 , tmp )
c98  format ('iam=',i3,' counter=',i5,' wait_t_time=',1p,e15.8)
call gsync()
if (trans_id .ne. nodes) then
    call error_stop(20,idummy,idummy,check_file,error_flg)
endif
call error_exit(error_flg)
c   if (iam.eq.0) then
c       write (*,99) trans_type,trans_id
c99  format('subroutine trans_check',i5,' end trans_id=',i3)
c   endif
c
return
end
```

```
C=====
c      subroutine wait_time
C=====
c
c      subroutine wait_time(w_time)
c        include 'fnx.h'
c        real*8 w_time,s_time,t_time
c        logical*4 wait_flg
c
c        wait_flg = .ture.
c        s_time = dclock()
c        do while ( wait_flg )
c          t_time = dclock() - s_time
c          if ( t_time .gt. w_time ) then
c            wait_flg = .false.
c          endif
c        enddo
c
c        return
c      end
```

付録3 メッセージパッシングの手法

Intel paragonでは、各々のプロセッサ（ノード）がメッセージパッシングとよばれる手法でデータを相互に交換することによって処理が進められる。メッセージパッシングの基本は、ノード間の一対一のメッセージの送受信である。ブロードキャストや入出力もこの組み合わせで実現されている。後者は、I/Oノードへのメッセージの送受信によって行われる。送受信は、送り側ノードと受け側ノードで一対の関数、またはサブルーチンを実行することによって行われる。

Intel paragonのnxライブラリで通常用いられるメッセージの送受信の手続きには、ブロッキング型通信のcsend、crecvサブルーチンとノンブロッキング型のisend、irecv関数の2系統があり、高速かつ信頼性の高い通信を行うには、これらの特性を良く理解することが重要である。

メッセージの送受信とは、送り側ノードのプログラム中の変数や配列の形でメモリに保持されているデータを、受け側ノードのメモリへ転送する手順である。（Fig.A3.1）のように送り側ノードをノード0、受け側ノードをノード1としてその間のcsend、crecvによる通信を考える。まず、送り側ノードでcsendを発行すると、送信データを保持する変数や配列の内容が、オペレーティングシステムが通信用に確保しているバッファ（メッセージバッファ）にコピーされる。送り側ノードはこのデータを約1kBのパケットに分けてネットワークを通して受け側ノードに向けて送信する。受け側ノードはネットワーク上のデータを受け取りメッセージバッファに格納し、完了した時点でそのことを知らせる完了ステータスを送り側ノードに返す。完了ステータスを受けるとcsendは通信が正常に終了したものと判断して終了するので、プログラムは次の処理に移る。受け側ノードのcrecvは、メッセージバッファの内容を変数や配列にコピーし、完了ステータスをcrecvに返す。これによりcrecvは通信が正常に終了したと判断して終了し、プログラムは次の処理に移る。

以上は、通信するデータのサイズがメッセージバッファのサイズよりも小さい場合である。より大きなデータの通信を行う場合の一般的な手順を（Fig.A3.2）に示す。メッセージバッファのサイズは、全体の大きさが決められたサイズをプログラムに属するノードひとつひとつに均等に割り付けるという方式で決められる。デフォルトのメッセージバッファのサイズは1MBであり、プログラムが10ノードで実行される時、ある一つの相手ノードとの通信のために割り当てられるサイズは100kB、100ノードの場合は10kBとなる。

データのサイズがメッセージバッファのサイズよりも大きい場合、送信用の変数や配列の内容は、メッセージバッファのサイズごとに分割されてコピーされる。データが受け側データのメッセージバッファに格納されてそれに対する完了ステータスが返されると、送信データの次の部分がメッセージバッファにコピーされる。受け側ノードでcrecvが発行されていれば、メッセージバッファに受信されたデータはすぐに受信用の変数や配列にコピーされ、受け側ノードは次のデータをメッセージバッファに受信することができる。しかし、crecvが発行されていなかったり、何らかの原因で受信したデータをすぐに処理出来ない場合、新しいデータを受信することは不可能である。完了ステータスが得られなかった場合、送り側ノードは通信のリトライを繰り返す。このようにして最後のメッセージバッファの内容が受け側ノードに転送された時点で、csendは通信が正常に終了したものと判断して終了し、プログラムは次の処理に移る。

ここに、プログラムを書く上で間違ひの原因となる問題がある。偏微分方程式を領域分割法で

解く場合には、タイムステップ毎に隣り合うノードがデータを交換しながら処理を進めることが多い。この時、

ノード0のプログラム	ノード1のプログラム
...	...
csend(type, A, count, 1, ptype)	csend(type, A, count, 0, ptype)
crecv(type, B, count)	crecv(type, B, count)
...	...

のようなプログラムは、原理的にハングアップする可能性がある。その理由は、フローチャートを (Fig.A3.3) に示すように、各々のノードで csend が完了しない限り crecv は発行されず、crecv が発行されない限り csend が完了しないからである。ただし、通信量がメッセージバッファのサイズよりも小さい時は、crecv が発行されていなくてもデータは全て受け側のメッセージバッファにコピーされて csend は完了してしまう。このプログラムは、転送データの量が十分小さい時は動作するが、メッセージバッファのサイズはノード数に反比例して減少することを考えると、少ないノード数では動作してもノード数を増加したときにデッドロックを発生すると考えられる。

この問題は、(Fig.A3.4)に示すように、通信の開始と完了の確認をわければ解決する。nxライブラリではcsendとcrecvの代わりにisend、irecv関数とmsgdone関数を用いて実現することができる。この場合、isendの完了を待たずにirecvを発行することができるから、ノード0はノード1へデータを送りながらノード1からのデータを受け取ることができ、デッドロックなしにデータをやり取りできる。送信用配列、変数の内容をメッセージバッファにコピーする作業、送り側ノードのメッセージバッファの内容を受け側ノードに転送する作業、受信メッセージバッファの内容を受信用配列、変数にコピーする作業は、プログラム本体とは関係無く（非同期に）実行される。なお、isend、irecv関数が発行されてから対応するmsgdone関数が完了するまでの間、送受信用配列、変数のデータは不定となり、この値を変更したりしてはいけない。

`isend`、`irecv`、`msgdone`関数を用いる方式では、通信の完了を判断する部分をプログラムに付け加えなければならずプログラムが複雑になる問題がある。同じ処理は、(Fig.A3.5) に示すように `isend`と`crecv`の組み合わせで実現することができる。(Fig.A3.4) の例では、受信が完了した場合には送信は完了しているので、両方の確認を行うことは冗長である。そこで、`isend`を発行したあとで`irecv`と`msgdone`の一組を連続して発行するとすれば、これは`crecv`を発行したのと同じことであるから、`crecv`で置き換えることができる。この方法により、メッセージバッファのサイズとデータのサイズに依存せずデッドロックすることなく通信を行うことができる。

ノード 0 のプログラムの実行位置 ... do i=1,imax ... 処理 ... (ループ $i+1$ 番目) isend(type, A, count, 1, ptype) ← crecv(type, B, count) end do	ノード 1 のプログラムの実行位置 ... do i=1,imax ... 処理 ... isend(type, A, count, 0, ptype) crecv(type, B, count) ← (i 番目) end do
---	---

しかし、もしこの処理が頻繁に繰り返される時は、問題が起こることがある。ノード0とノード1のプログラムの実行の順序は陽に同期を取らない限り保障されないので、ノード0でのプログラムの方がより早く進行して、上の例のようにノード0はすでに*i+1*番目のループの処理を実行してisend(type, A, count, 1, ptype)を実行しようとするのに、ノード1では*i*番目のループのcrecv(type, B, count)が完了していない状態が起こるかも知れない。ノード0のループが*i+1*番目へ進むには、ノード1のisend(type, A, count, 0, ptype)とノード0のcrecv(type, B, count)が完了していれば良いからである。一方、これはノード0の*i*番目のループのisend(type, A, count, 1, ptype)が完了していないことを意味するので、さらにノード0でisend(type, A, count, 1, ptype)を発行しようとすると通信用配列、変数のデータが破壊され、プログラムの動作は異常になる。

isend、irecv、msgdoneの組み合わせを用いて通信すればこの問題は起こらない。また、ループを進める度にノード間の同期を取れることによっても回避できる。ノード0（1）が発行するisendのメッセージタイプをループの回数毎に変えられることでも回避できるが、極端に多くのisendが発行された場合の動作は、メッセージタイプが重複する可能性があることや、一度にシステムが保持できる通信要求の数に制限があるために確実とは言えない。

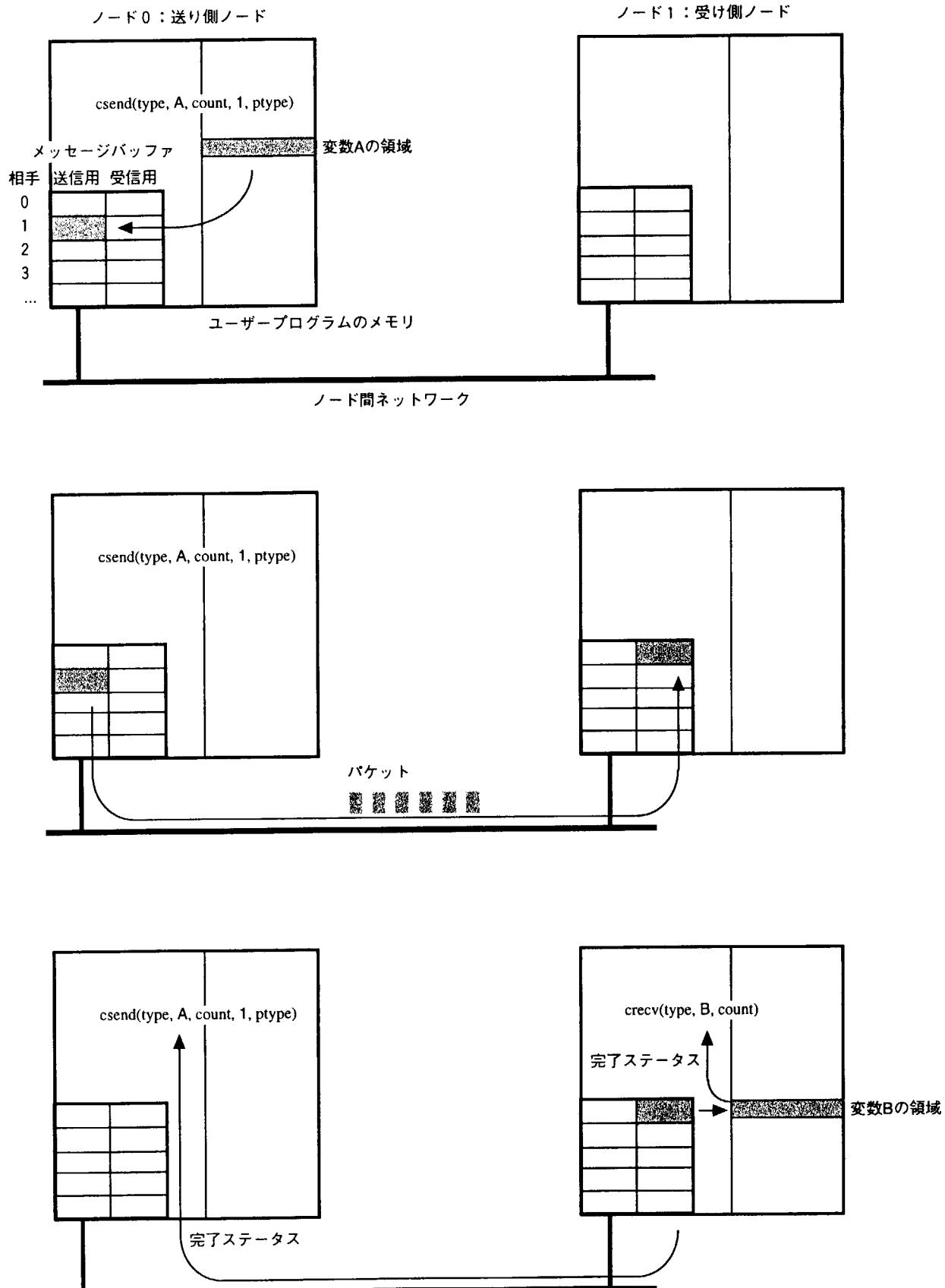
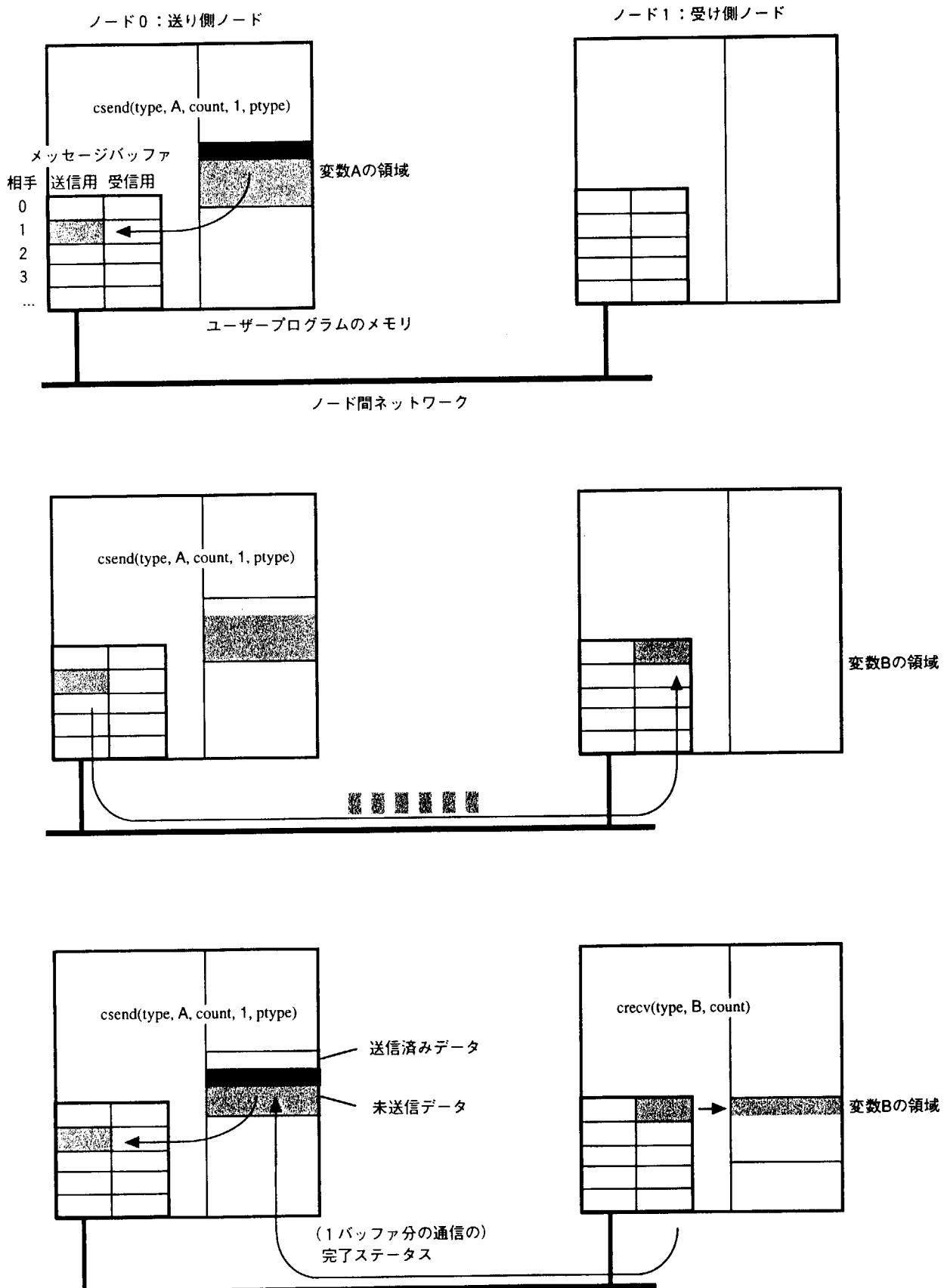


Fig.A3.1 ノード間のデータの送受信



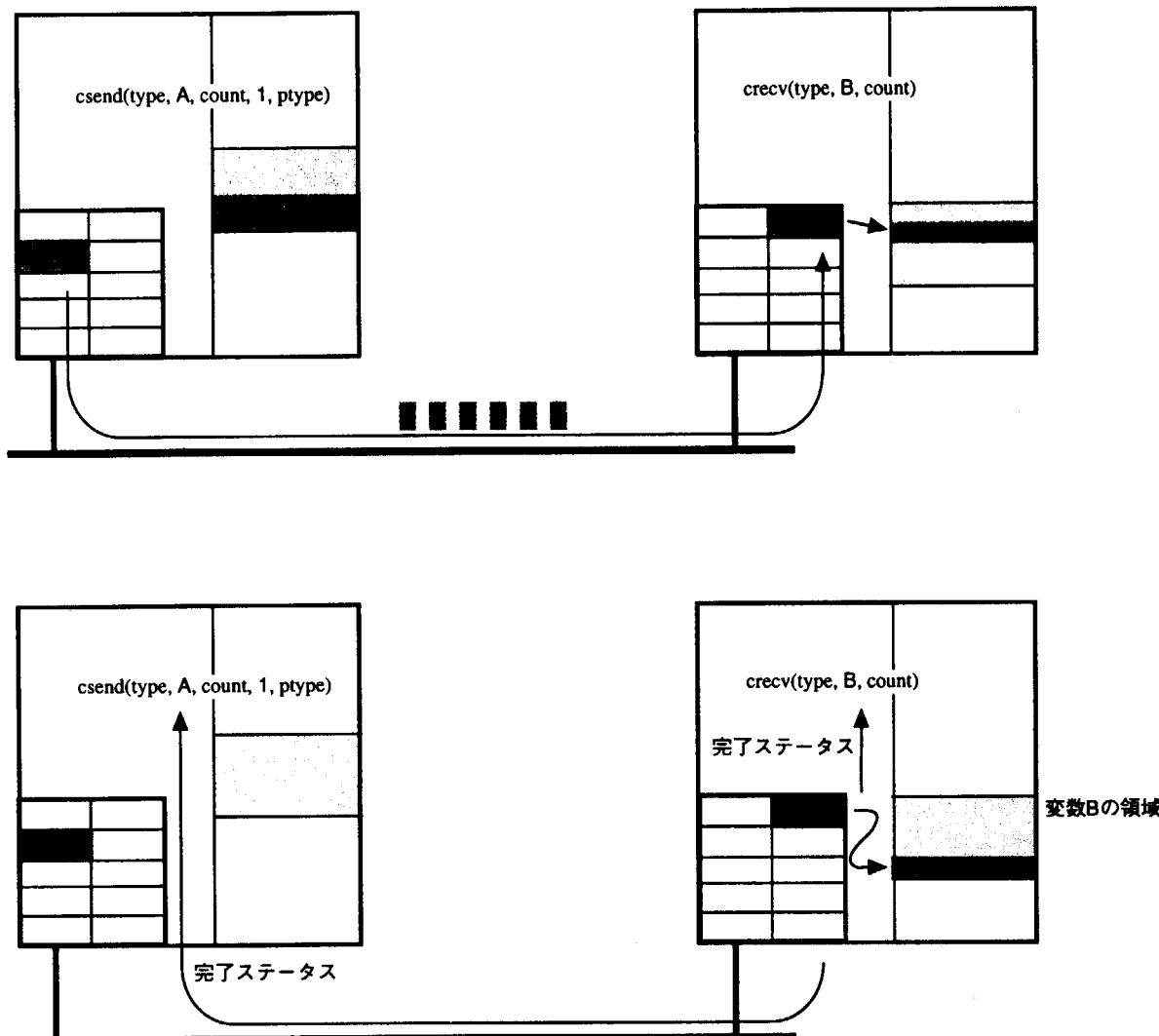
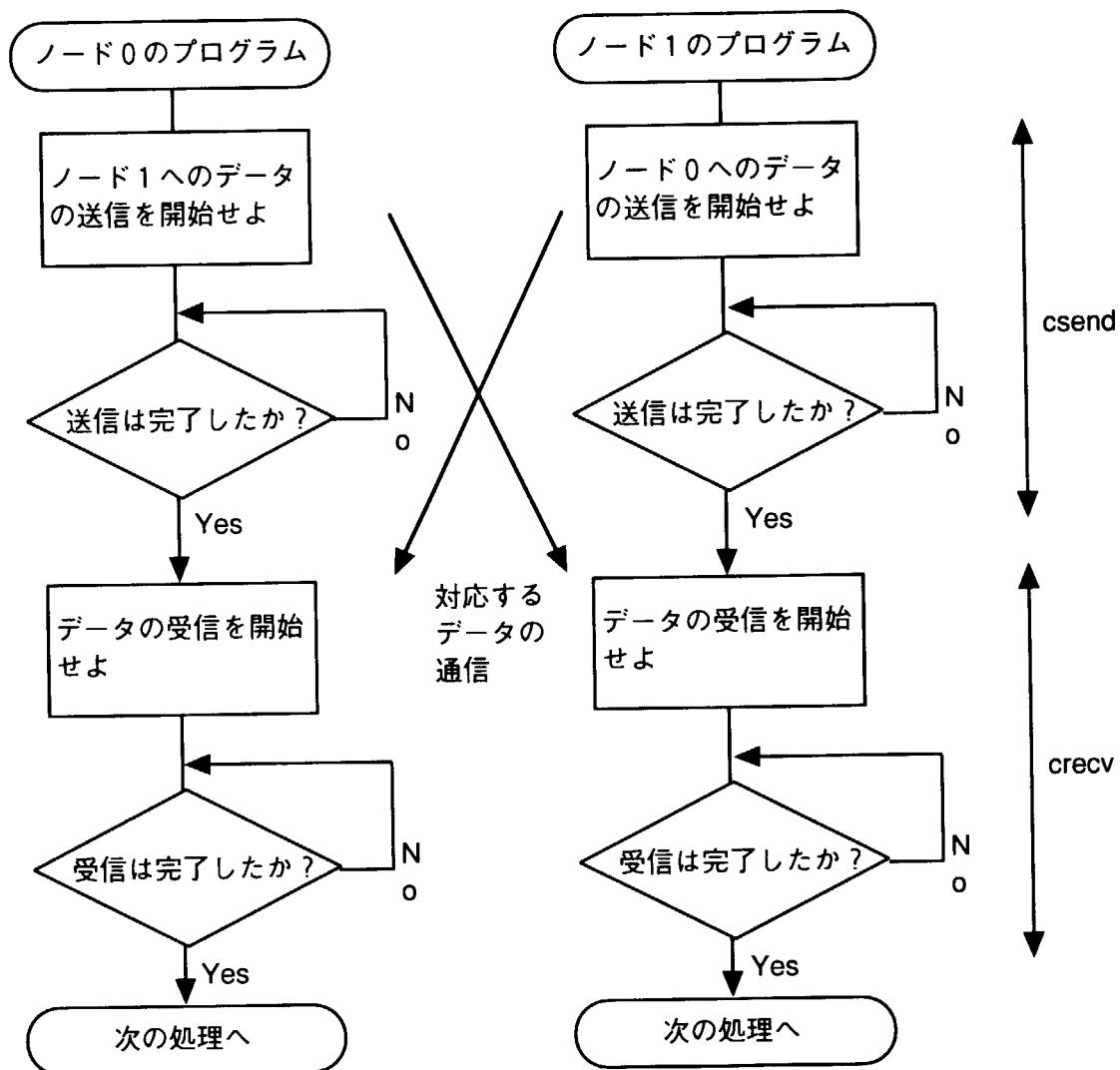


Fig.A3.2 データ量が多い場合のノード間の送受信

FigA.3.3 `csend`と`crecv`による通信

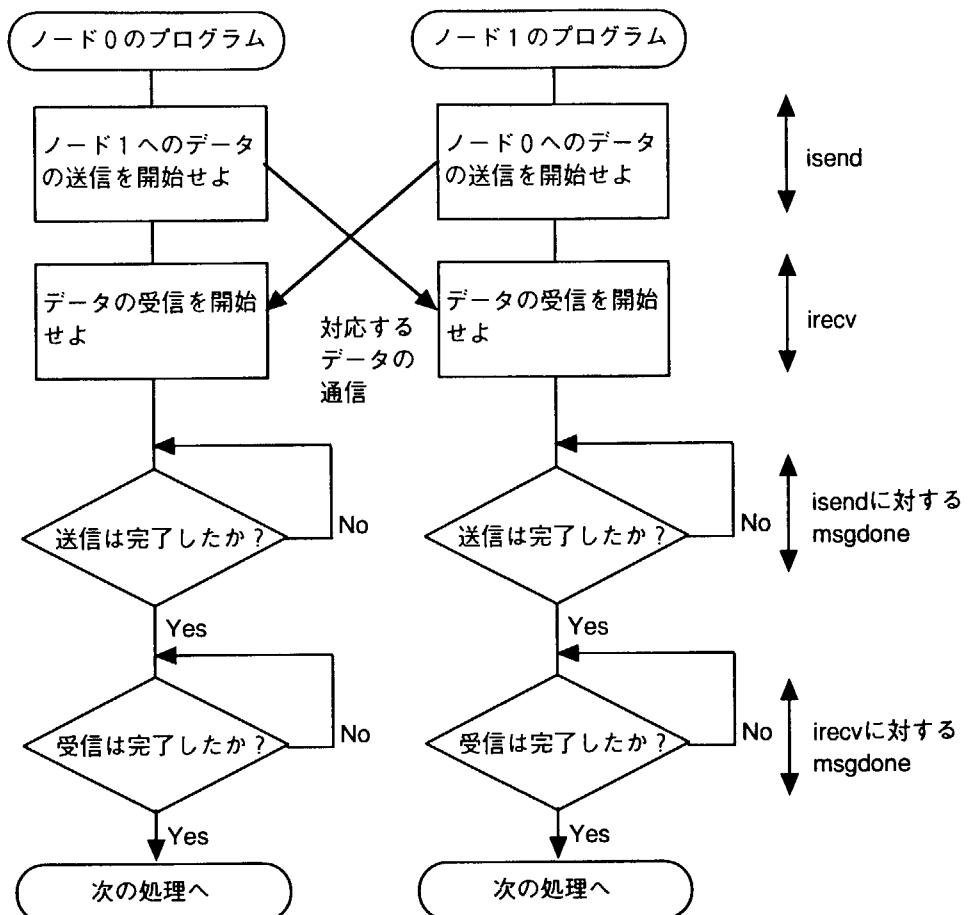


Fig.A3.4 isend、irecv、msgdoneによる通信

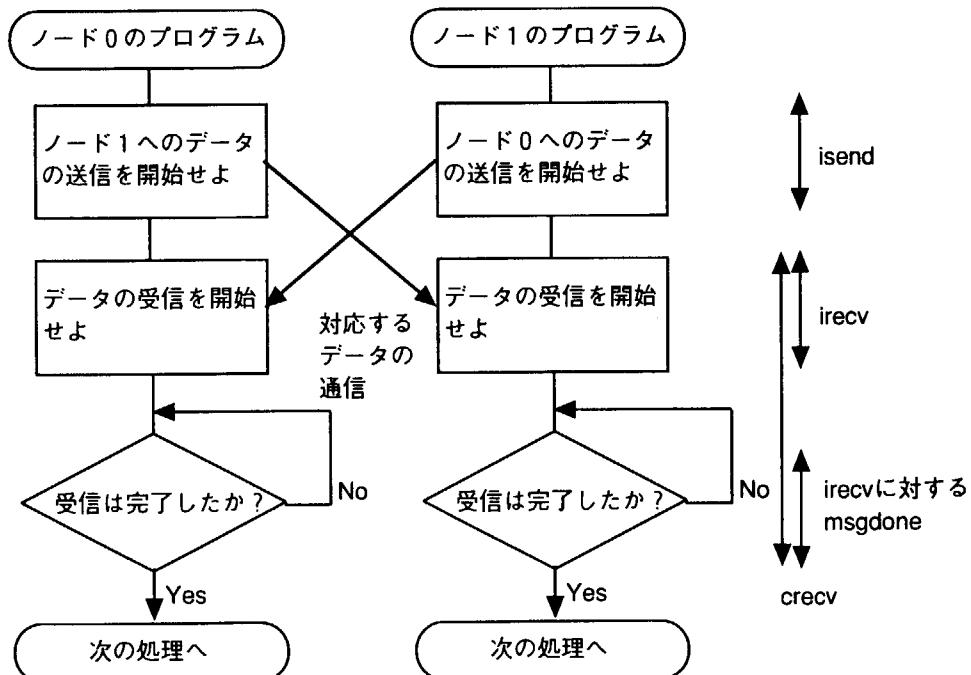


Fig.A3.5 isend、crecvによる通信

付録4 並列入出力の方法

分散記憶型超並列計算機では、計算用ノードの数に比例して保存しなければならないデータ量が非常に多いのに対して、ディスク等のI/Oデバイスの数は限られていて、しかもノード間のネットワークやローカルエリアネットワークを介して接続されているので、入出力が計算機の性能を制約する大きな要因になりうる。

fortranによる科学技術計算を行う時、データの入出力はシーケンシャルにできれば良い場合が大半と考えられる。(Fig.A4.1)に示すように、各々のwrite文ごとにそれぞれ異なるサイズのデータが順々にファイルに書き込まれる。計算に対してディスクへの書き込みの速度は通常遅いので、書き込みの完了を待つ時間が処理速度を制約する原因になる。

これを解決するため、1台のコンピュータに複数台のディスクを装備することがある。

(Fig.A4.2)に示すように、データは各ディスクに分割して書き込まれる。ユーザーからは、このようなデータファイルはあくまでも単一のファイルに見えるような技術が使われていて、物理的にどのディスクに書き込まれるかはデータのサイズ等から自動的に計算され、プログラムからは見えないのが普通である。これは一種の出力の並列化であり、書き込みの時間がディスクの台数に反比例して減少し高速化されると期待される。

超並列計算機で計算したデータを出力する場合は新たな問題が発生する。Intel Paragonでは、ディスクは専用のI/O計算ノードに接続されているので、計算ノードからメッセージパッキングによる通信を行うことではじめて読み書きが可能となる。この場合、より多量のデータを通常のメモリ、MPUとディスクの間の転送速度よりも遅いノード間のネットワークを介して出力しなければならない。さらに、複数の計算ノードからひとつのファイルに読み書きを行うために、アクセスの調停を行わなければならない。

これを解決するため、Intel paragonでは、複数の計算用ノードから複数のI/Oノードとの間で読み書きを行うためのpfsが実装されており、メッセージパッキングを応用したデータの入出力の方法が定義されている。pfsに対する出力を開始するには、gopenサブルーチンによりファイルをオープンし、その後cwriteサブルーチンまたはiwrite関数によって書き込みを行う。gopenサブルーチンは全てのノードで実行しなければならないグローバル命令である。

さて、(Fig.A4.3)に示すように、各々のノードが各々のwrite文ごとに異なるサイズのデータの出力をを行うとすると、ノード1はノード0の書き込みが終了しなければディスクファイル上の何処に書き込みを行って良いかわからない。gopenサブルーチンでM_SYNCモードを選択すると、アクセスの調停が行われ、各ノードからの読み書きが自動的に順番に行われるようになる。すなわち、(Fig.A4.3)に示すようにノード0の出力が終わった時点で、現在のディスクファイル上の書き込み位置(ファイルポインタの値)がノード1に送られ、ノード1からの出力が行われ、ファイル上には発行したwrite文毎、ノード順にデータが整列して保存される。しかしながら、この方法では同時に1台のノードしか読み書きを行うことが出来ない。ノード1はノード0の書き込みが完了しないと出力を行うことができない。他の全てのノードが読み書きを完了するまで次の処理に進めないので並列処理のメリットを發揮できない。

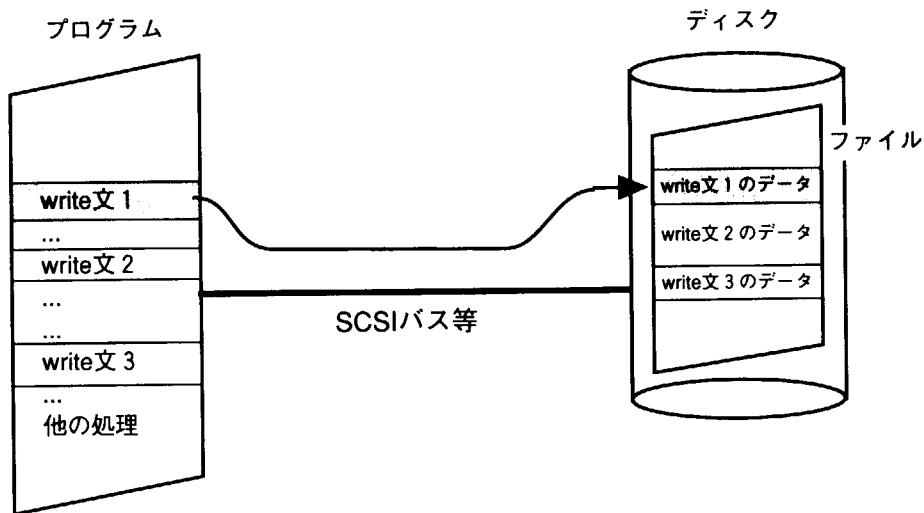
各ノードが発行するwrite文が厳密に同じサイズのデータの出力をを行うことが仮定できる場合は、M_RECORDモードを用いてより効率的な出力をを行うことができる。この場合、各ノードは他

のノードからの情報なしにそれぞれのwrite文におけるディスクファイル上の書き込み位置（ファイルポインタの値）を計算することができる。したがって、(Fig.A4.4) に示すように各ノードは独立に書き込みを行うことができる。図のように書き込みの進捗がノード毎に異なっていても構わない。出力ファイルはこれまでの場合と異なってシーケンシャルにアクセスされない。I/Oノードは自動的にファイルポインタを移動させて所定の場所にデータを書き込む処理を行う。

M_RECORDモードでの出力の対象となるディスクファイルは、物理的に (Fig.A4.2) のような分割されたディスクのどちらかでも良いし、(Fig.A4.5) のように別のI/Oノードに接続されているディスクでも良いと考えられる。pfsはこのような処理を実現するためのもので、複数のI/Oノードに接続されているディスクに分割されているファイルを、ユーザーからは論理的には一続きのシーケンシャルファイルのように見えるようにする。pfsを活用すると、ノード0からの出力はI/Oノード0で、ノード1からの出力はI/Oノード1で、独立に処理される。このように効率的に並列処理する方法をスライティングと言う。

pfsを用い並列出力の効率を高めるための注意点は、スライティングがサイズ64kB毎に機械的に行われることである。各ノードからのひとつのwrite文で出力されるデータサイズが64kB（以下）の場合は、各ノードからのデータが整然と各I/Oノードに送られ処理されるが、64kBを超える場合には別の出力ノードに分割されて振り分けられてしまう。結局、(Fig.A4.6) の例ではそれぞれのwrite文毎にデータがノード0と1の両方に送られ、ノード間ネットワークや各々のノードからの出力要求を受けるI/Oノードの負荷が増加する。これは処理速度の低下とプログラムの信頼性の低下を引き起こす可能性がある。スライティングされて実際に出力が行われる先のI/Oノードをプログラムから指定することはできない。各write文のデータサイズを厳密に64kBに一致させることは容易ではないので、出力データを64kBずつバッファリングするサブルーチンが必要かも知れない。

以上のようなcwrite、iwriteサブルーチン、関数を用いて出力されたデータは、出力文と完全に一対一対応するcread、ireadサブルーチン、関数で読み込まなければならない。並列出力はI/Oノードへのメッセージパッシングであるとともに、ファイルを介したプログラムからプログラムへのメッセージパッシングを考えることもできる。



(Fig.A4.1) ディスクへのシーケンシャルファイルの書き出し

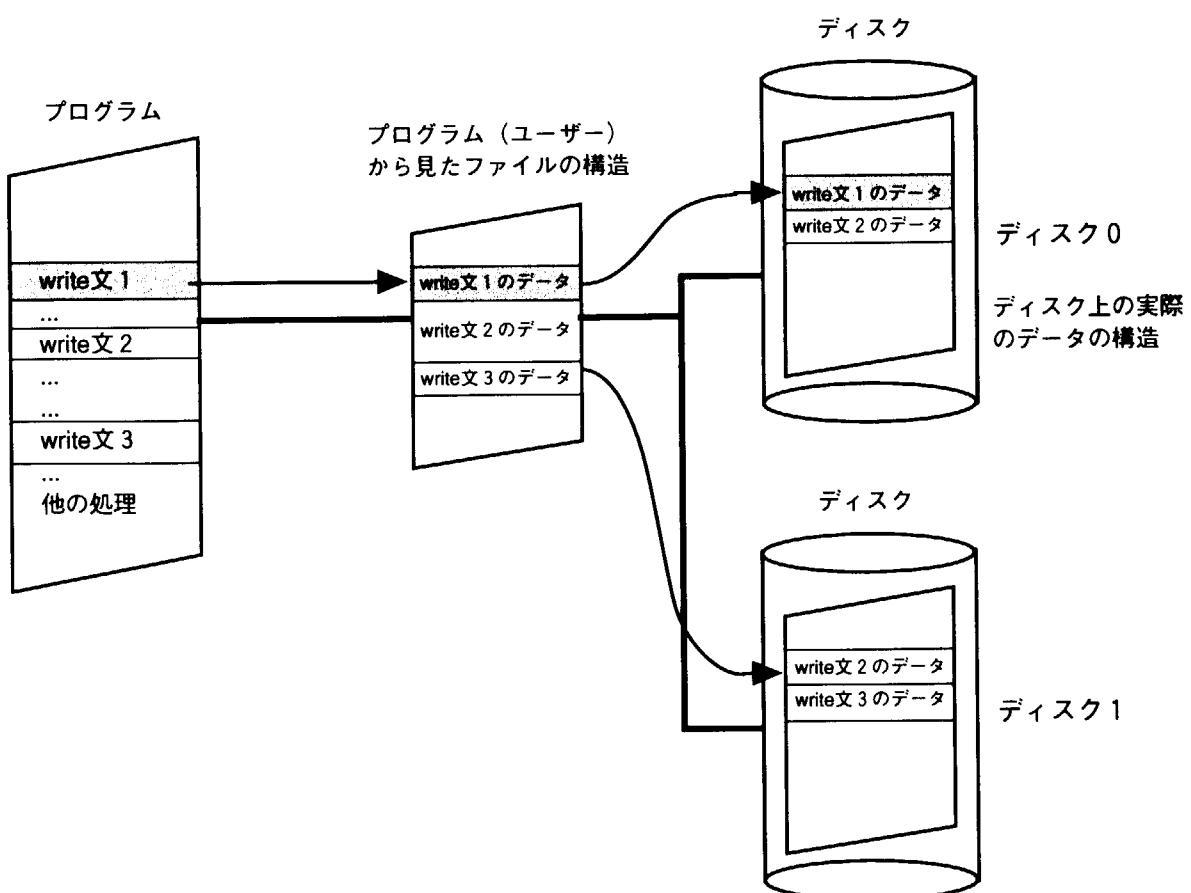


Fig.A4.2 複数ディスクへのシーケンシャルファイルの書き出し

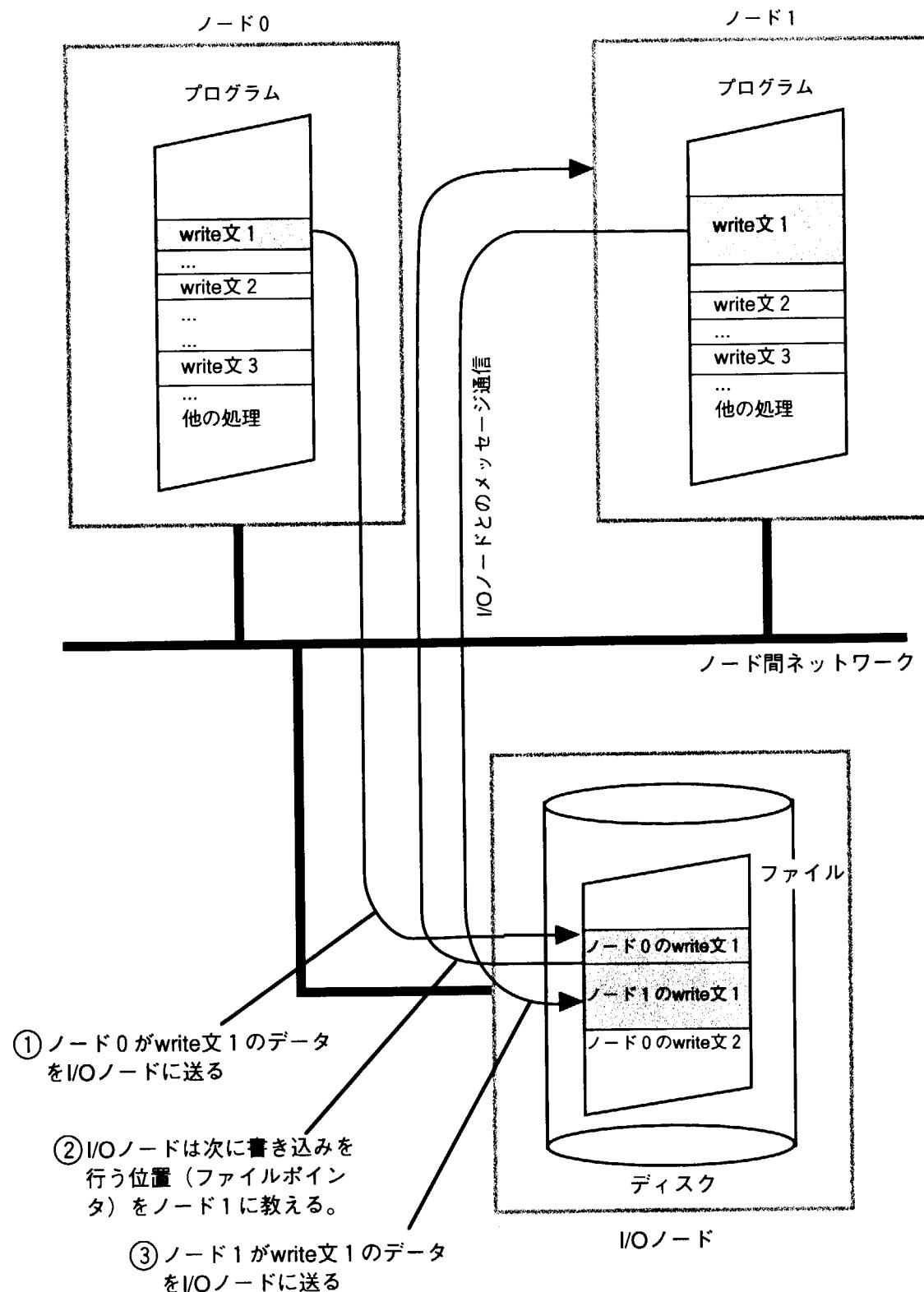
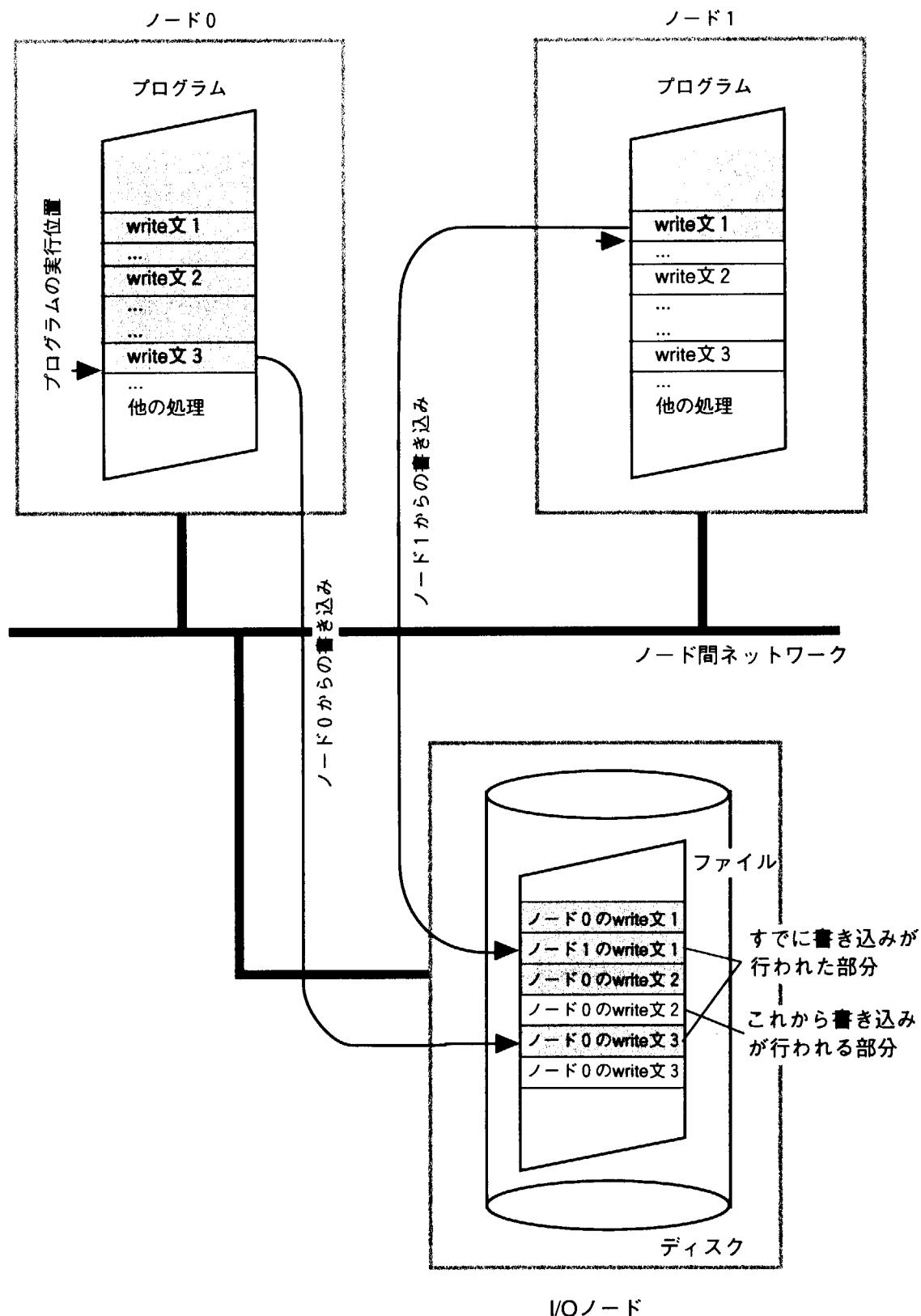
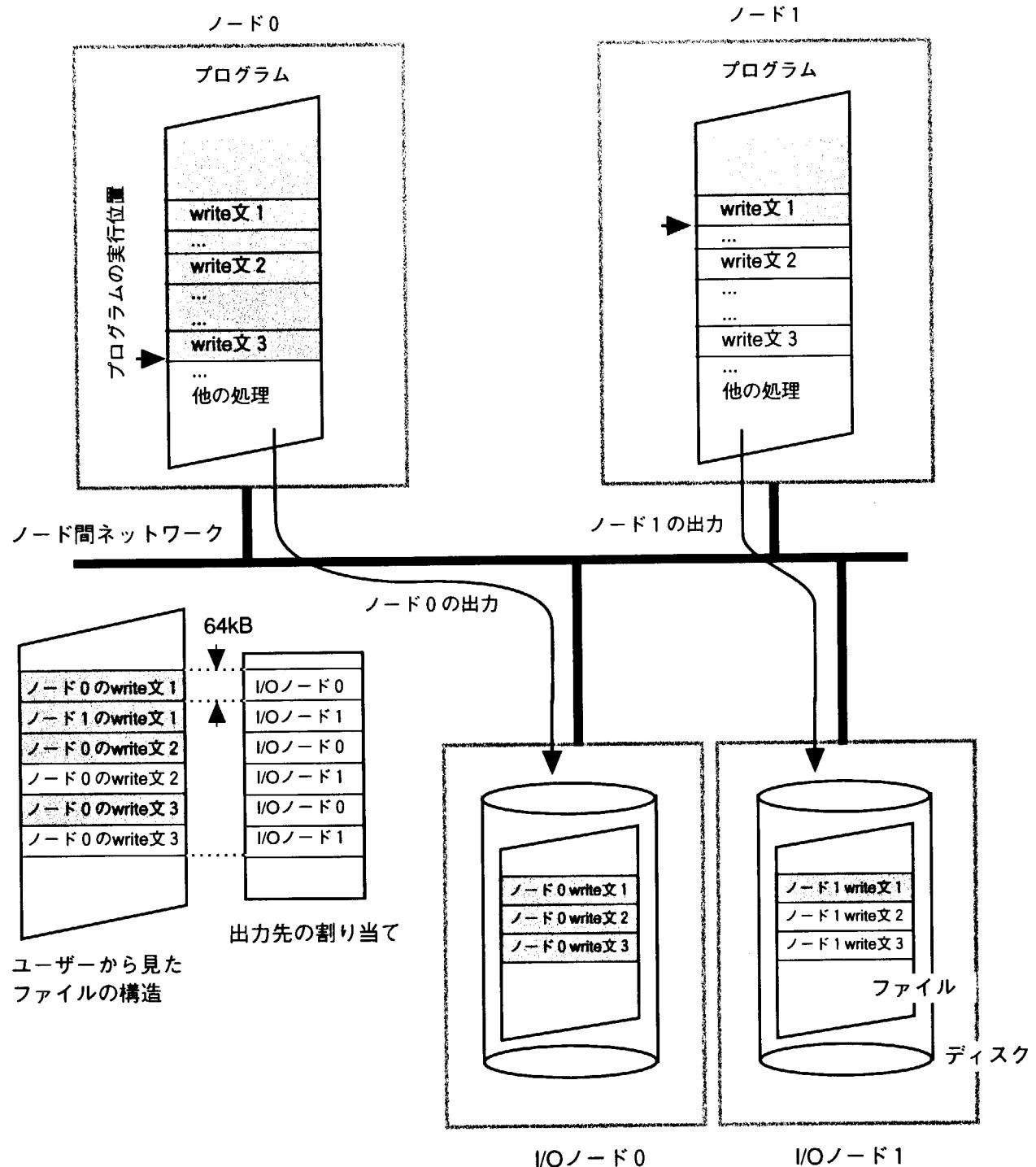


Fig.A4.3 M_SYNC I/Oモードによる複数ノードからディスクへの書き出し



FigA4.4 M_RECORD I/Oモードによる複数ノードからディスクへの書き出し



FigA4.5 pfsを用いた複数ノードから複数I/Oノードへの書き出し

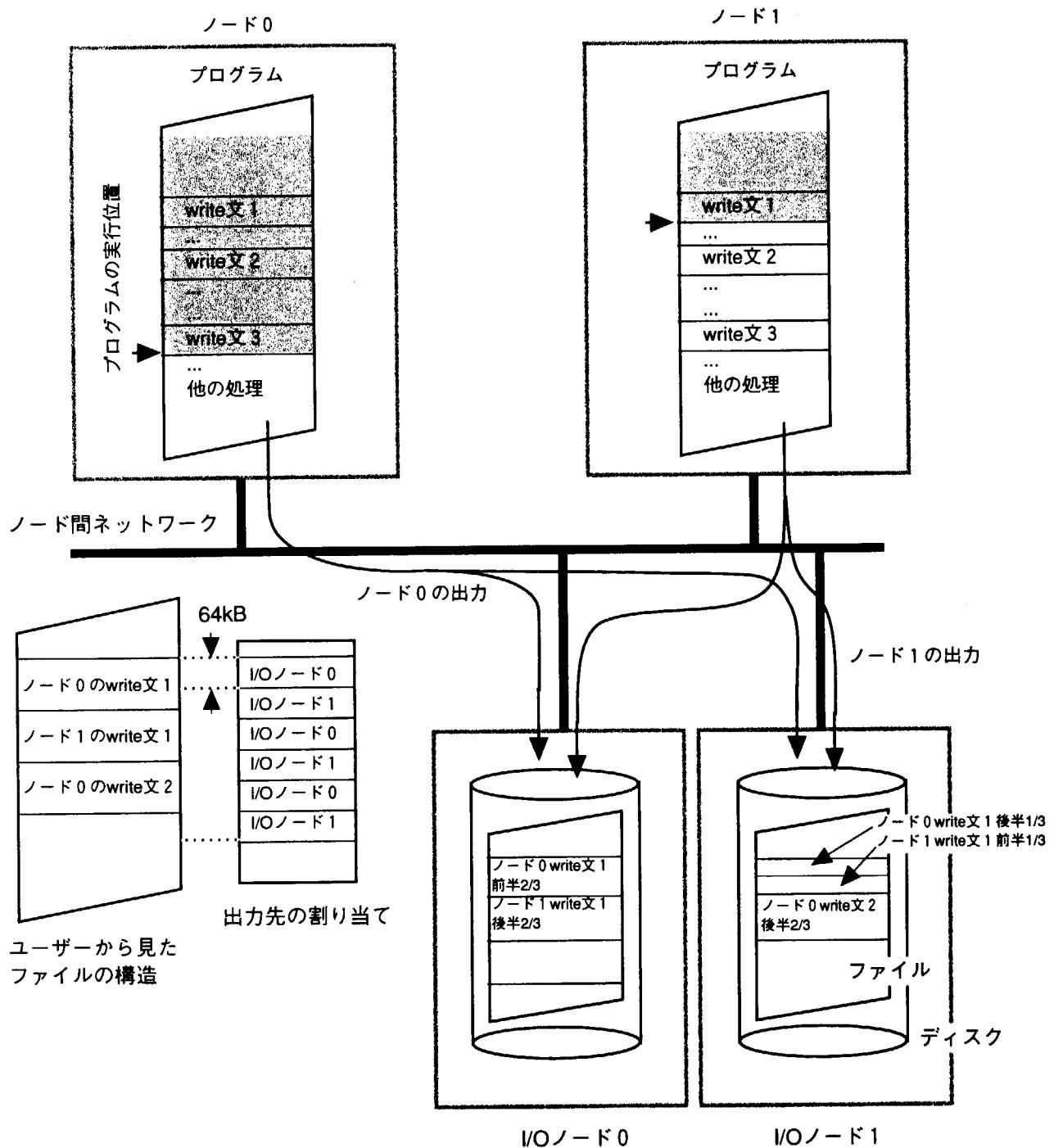


Fig.A4.6 スライトサイズと出力サイズの不一致時の処理

付録5 割り込み処理によるプログラム実行の制御

プログラムを実行中、入力データの異常、計算結果の異常、入出力装置の異常によりこれを停止しなければならないことがしばしばある。しかし、並列計算機のひとつノードで実行されているプログラムに異常が発生した時、ジョブを安全に停止させるには、すべてのノードで実行されているプログラムを停止させることが必要となる。このような技術は自明ではない。

メッセージパッシングによる通信は、送り側ノードと受け側ノードが一対の`csend`、`isend`と`crecv`、`irecv`を発行することによって成立する。すなわち、各ノードが`csend`、`isend` (`crecv`、`irecv`) を発行する時、相手側ノード上でプログラムが正常に動作して`crecv`、`irecv` (`csend`、`isend`) が発行されることが前提である。このため、通常のプログラムでは、(Fig.A5.1) のように他のノードで異常が発生していることを積極的に知ることは不可能である。あるノードが`stop`文を実行したり、実行時エラーで終了しても、他のノードはそのことが分からないのでそのまま実行を続ける。その後のプログラム中に通信が含まれていれば、これは決して完了しないので、プログラムは確実にハングする。異常が発生したノードからこのことを他のノードに伝えようとしても、受け側のノードで適当な`crecv`、`irecv`が発行されていなければ内容は伝わることがない。

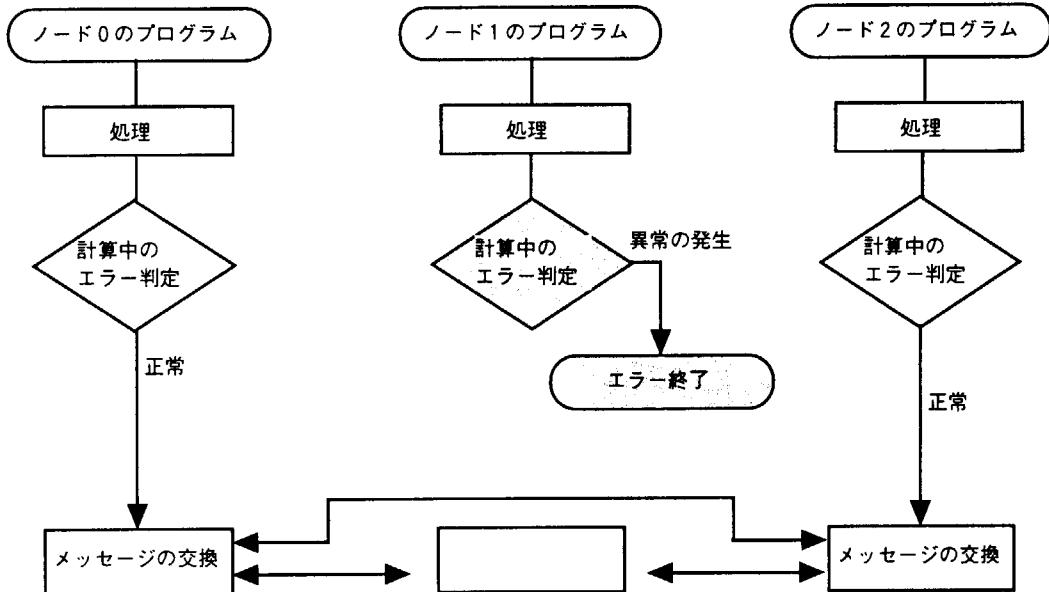
このための対策の一つは、プログラム中の適当な場所で全てのノードがエラー情報を交換しあうことである (Fig.A5.2)。これによって、あるノードで発生したエラーを全てのノードが察知してジョブを終了させることができる。しかしプログラム中にこのための処理を多数回加えると、プログラムが複雑になる上、エラー情報の交換に要する通信時間のために処理速度が損なわれる。特にすでに停止てしまっているノードがあった場合、これを判定する唯一の方法は、所定時間応答がなかった場合にエラーと判定する方法だけである。このための待ち時間は処理速度を損なうと予想される。

Intel Paragonのnxライブラリでは、割り込み型のメッセージパッシングの手続きが定義されていて、これを用いるとより柔軟なエラー処理が行えると考えられる。割り込み処理とは、ハードウェアの信号やソフトウェアのある特定の命令（割り込み）を用いて、現在実行されているプログラムを一時中断させて別のプログラムを実行することである。

割り込みが発生した時に実行されるプログラムを割り込みハンドラと呼ぶ。Intel Paragonでは送り側ノードが割り込み型のメッセージを送ると、受け側ノードでどのような処理を行っているかに関わらず割り込みハンドラを起動されるようにすることができる。プログラム本体は、エラーを検出した任意の時点で割り込み型のメッセージを発行するだけで良く、割り込みハンドラはプログラムを終了させる処理のみを行えば良い (Fig.A5.3)。この方法ではエラー情報の交換のための待ち時間もないので高速な処理が行えると考えられる。

割り込みを活用するとより多様な処理を行うことができる。例えば、割り込みハンドラにその時点のデータを書き出す処理を加えると、エラーの原因の解析に便利である。長時間のジョブの場合、書き出されたデータをもとにリスタートできればCPU時間を有効に活用することもできるだろう。反面、複雑なプログラムには、メモリアロケーションやpfsへの読み書きなど中断が許されない処理（クリティカルセクション）が存在し、このような状況下で割り込み処理を活用するにはより高度なプログラム技術が必要になる。

MPIを用いる場合は、エラーが発生したノードで`MPI_ABORT`関数を呼び出すことにより、単純にすべてのノードで実行されているプログラムを停止させ、ジョブを終了させることができる。



(Fig.A5.1) 一つのノードでのエラーによるプログラムハンギングの例

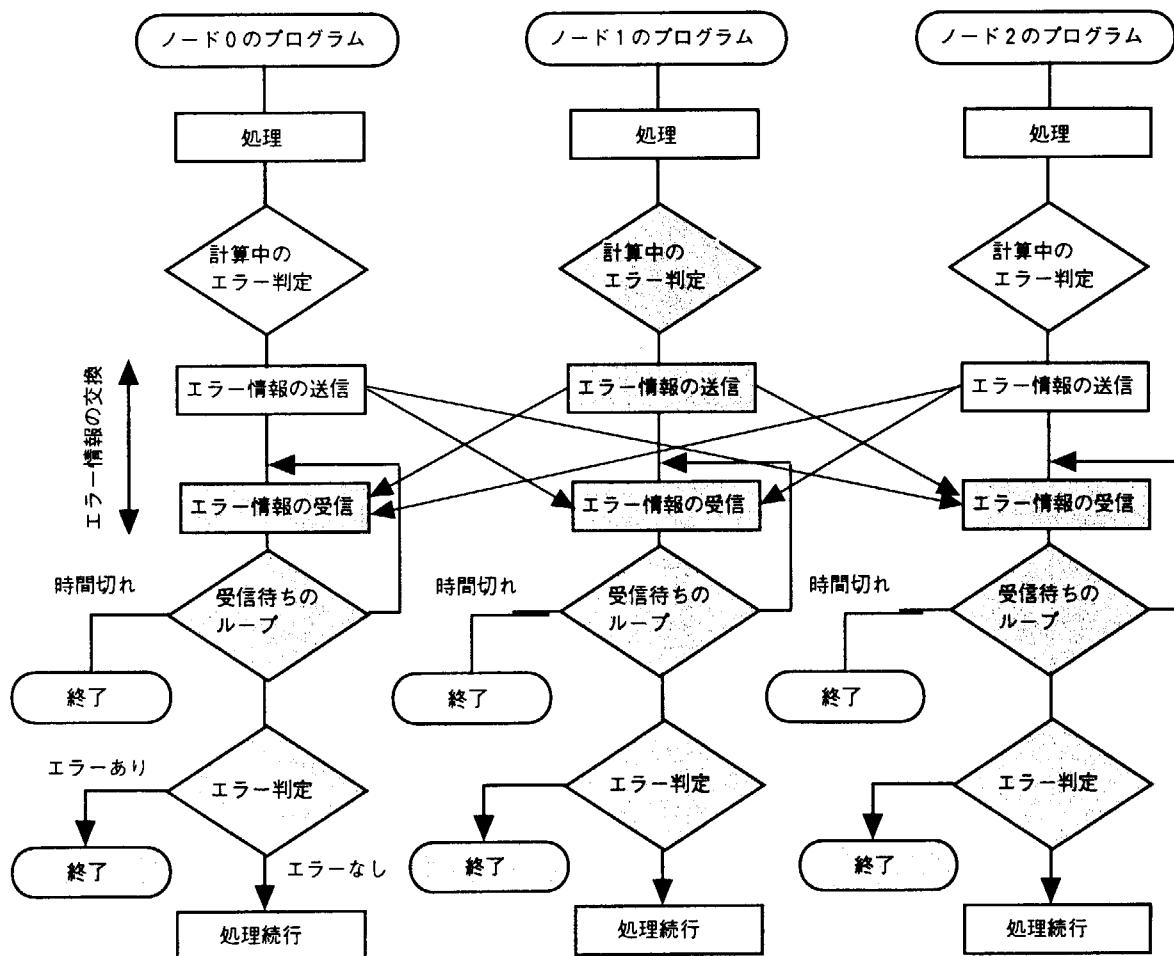


Fig.A5.2 エラー情報の交換を加えたプログラム

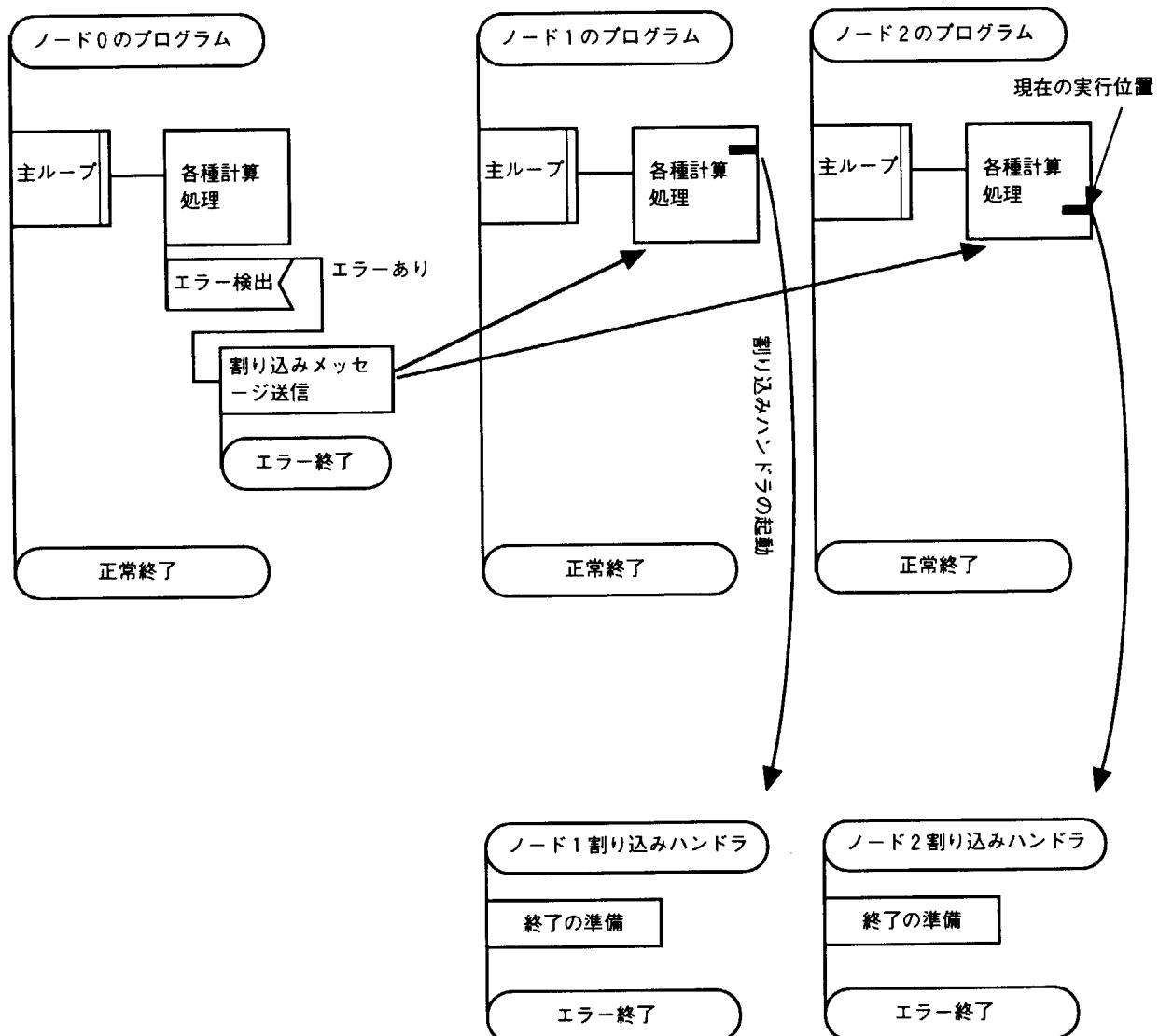


Fig.A5.3 割り込みによるエラー処理

This is a blank page.

国際単位系(SI)と換算表

表1 SI基本単位および補助単位

量	名称	記号
長さ	メートル	m
質量	キログラム	kg
時間	秒	s
電流	アンペア	A
熱力学温度	ケルビン	K
物質量	モル	mol
光度	カンデラ	cd
平面角	ラジアン	rad
立体角	ステラジアン	sr

表3 固有の名称をもつSI組立単位

量	名称	記号	他のSI単位による表現
周波数	ヘルツ	Hz	s ⁻¹
力	ニュートン	N	m·kg/s ²
圧力、応力	パスカル	Pa	N/m ²
エネルギー、仕事、熱量	ジュール	J	N·m
功率、放射束	ワット	W	J/s
電気量、電荷	クーロン	C	A·s
電位、電圧、起電力	ボルト	V	W/A
静電容量	ファラード	F	C/V
電気抵抗	オーム	Ω	V/A
コンダクタンス	ジーメンス	S	A/V
磁束	ウェーバ	Wb	V·s
磁束密度	テスラ	T	Wb/m ²
インダクタンス	ヘンリー	H	Wb/A
セルシウス温度	セルシウス度	°C	
光束度	ルーメン	lm	cd·sr
照度	ルクス	lx	lm/m ²
放射能	ベクレル	Bq	s ⁻¹
吸収線量	グレイ	Gy	J/kg
当量	シーベルト	Sv	J/kg

表2 SIと併用される単位

名称	記号
分、時、日	min, h, d
度、分、秒	°, ', "
リットル	l, L
ト	t
電子ボルト	eV
原子質量単位	u

$$1 \text{ eV} = 1.60218 \times 10^{-19} \text{ J}$$

$$1 \text{ u} = 1.66054 \times 10^{-27} \text{ kg}$$

表4 SIと共に暫定的に維持される単位

名称	記号
オングストローム	Å
バーン	b
バール	bar
ガル	Gal
キュリ	Ci
レントゲン	R
ラド	rad
レム	rem

$$1 \text{ Å} = 0.1 \text{ nm} = 10^{-10} \text{ m}$$

$$1 \text{ b} = 100 \text{ fm}^2 = 10^{-28} \text{ m}^2$$

$$1 \text{ bar} = 0.1 \text{ MPa} = 10^5 \text{ Pa}$$

$$1 \text{ Gal} = 1 \text{ cm/s}^2 = 10^{-2} \text{ m/s}^2$$

$$1 \text{ Ci} = 3.7 \times 10^{10} \text{ Bq}$$

$$1 \text{ R} = 2.58 \times 10^{-4} \text{ C/kg}$$

$$1 \text{ rad} = 1 \text{ cGy} = 10^{-2} \text{ Gy}$$

$$1 \text{ rem} = 1 \text{ cSv} = 10^{-2} \text{ Sv}$$

表5 SI接頭語

倍数	接頭語	記号
10 ¹⁸	エクサ	E
10 ¹⁵	ペタ	P
10 ¹²	テラ	T
10 ⁹	ギガ	G
10 ⁶	メガ	M
10 ³	キロ	k
10 ²	ヘクト	h
10 ¹	デカ	da
10 ⁻¹	デシ	d
10 ⁻²	センチ	c
10 ⁻³	ミリ	m
10 ⁻⁶	マイクロ	μ
10 ⁻⁹	ナノ	n
10 ⁻¹²	ピコ	p
10 ⁻¹⁵	フェムト	f
10 ⁻¹⁸	アト	a

(注)

- 表1～5は「国際単位系」第5版、国際度量衡局1985年刊行による。ただし、1eVおよび1uの値はCODATAの1986年推奨値によった。
- 表4には海里、ノット、アール、ヘクタールも含まれているが日常の単位なのでここでは省略した。
- barは、JISでは流体の圧力を表わす場合に限り表2のカテゴリーに分類されている。
- EC閣僚理事会指令ではbar、barnおよび「血圧の単位」mmHgを表2のカテゴリーに入れている。

換算表

力	N(=10 ⁵ dyn)	kgf	lbf
	1	0.101972	0.224809
9.80665	1	2.20462	
4.44822	0.453592	1	

粘度 1 Pa·s(N·s/m²) = 10 P(ボアズ)(g/(cm·s))

動粘度 1 m²/s = 10⁴ St(ストークス)(cm²/s)

圧力	MPa(=10 bar)	kgf/cm ²	atm	mmHg(Torr)	lbf/in ² (psi)
力	1	10.1972	9.86923	7.50062 × 10 ³	145.038
	0.0980665	1	0.967841	735.559	14.2233
	0.101325	1.03323	1	760	14.6959
	1.33322 × 10 ⁻⁴	1.35951 × 10 ⁻³	1.31579 × 10 ⁻³	1	1.93368 × 10 ⁻²
	6.89476 × 10 ⁻³	7.03070 × 10 ⁻²	6.80460 × 10 ⁻²	51.7149	1

エネルギー・仕事・熱量	J(=10 ⁷ erg)	kgf·m	kW·h	cal(計量法)	Btu	ft · lbf	eV	1 cal = 4.18605 J(計量法)
	1	0.101972	2.77778 × 10 ⁻⁷	0.238889	9.47813 × 10 ⁻⁴	0.737562	6.24150 × 10 ¹⁸	= 4.184 J(熱化学)
9.80665	1	2.72407 × 10 ⁻⁶	2.34270	9.29487 × 10 ⁻³	7.23301	6.12082 × 10 ¹⁹		= 4.1855 J(15 °C)
3.6 × 10 ⁶	3.67098 × 10 ⁵	1	8.59999 × 10 ⁵	3412.13	2.65522 × 10 ⁶	2.24694 × 10 ²⁵		= 4.1868 J(国際蒸気表)
4.18605	0.426858	1.16279 × 10 ⁻⁶	1	3.96759 × 10 ⁻³	3.08747	2.61272 × 10 ¹⁹		仕事率 1 PS(仏馬力)
1055.06	107.586	2.93072 × 10 ⁻⁴	252.042	1	778.172	6.58515 × 10 ²¹		= 75 kgf·m/s
1.35582	0.138255	3.76616 × 10 ⁻⁷	0.323890	1.28506 × 10 ⁻³	1	8.46233 × 10 ¹⁸		= 735.499 W
1.60218 × 10 ⁻¹⁹	1.63377 × 10 ⁻²⁰	4.45050 × 10 ⁻²⁶	3.82743 × 10 ⁻²⁰	1.51857 × 10 ⁻²²	1.18171 × 10 ⁻¹⁹	1		

放射能	Bq	Ci
	1	2.70270 × 10 ⁻¹¹
3.7 × 10 ¹⁰	1	

吸収線量	Gy	rad
	1	100
0.01	1	

照射線量	C/kg	R
	1	3876
2.58 × 10 ⁻⁴	1	

線量当量	Sv	rem
	1	100
0.01	1	

(86年12月26日現在)

パナソニックのスカラ超並列プロセッサ開発ガイド