

並列プログラミング入門 MPI版

2012年6月1日版

理化学研究所 情報基盤センター

青山 幸也

(このページは空白です。)



## ● はじめに ●

昔は並列計算機と言えば、特別な人たちが特殊な用途で使用する珍しいマシンであるという印象がありました。ところがここ数年、並列計算機の導入が急激に増加してきています。並列計算機の主要な用途として、科学技術計算とデータベース検索があります。このうち科学技術計算ではプログラムを並列化する必要があります。並列コンピュータに関する数多くの本が出版されていますが、並列コンピュータを作成する人のための並列アーキテクチャの本、並列コンパイラを作成する人のための本、情報系の大学院レベルを対象とした『理論』としての並列化の本などが多く、プログラムを並列化する方法について説明した本はあまり見当たりません。

並列計算機に関するこのような現状から、並列化を行う現場で役立つ分かり易いテキストを目標として本書を作成しました。しかし、並列プログラミングの技術はまだまだ発展途上であり、今後も多くの方が出てくると思われますので、本書もそれに合わせて内容を加筆修正していきたいと思います。

なお、本書の内容に関する責任は負いかねますのでご了承願います。

1995年 青山 幸也

本書は分散メモリ一型並列計算機上でプログラムを並列化する方法についてのガイドで、初版を作成してから早いもので約8年経ちました。初版は筆者がコンピュータ会社に勤務していたときにユーザー教育用として作成しました。私事で恐縮ですが、2002年秋に転職したのを機に、書名を改め、特定のマシンに固有の部分を除く、その他の部分も加筆修正する事にしました。このような経緯のため、内容に一部整合性がとれていない部分があるかと思いますが、今後少しずつ改善していきたいと思います。

2003年 青山 幸也

当資料は、日本アイ・ピー・エム株式会社の許諾を得て、「(虎の巻シリーズ2)並列プログラミング虎の巻 MPI版」を利用し、理化学研究所の責任で作成したものです。

本書は不定期に加筆修正しており、最新版は <http://acc.riken.jp/HPC/training.html> にあります。

# 目次

はじめに	i
第1章 並列プログラミングとは	1-1
1-1 分散メモリー型並列計算機とは	1-2
1-2 並列プログラミング・モデル	1-4
第2章 並列ジョブの実行方法	2-1
2-1 並列ジョブの実行方法	2-2
第3章 メッセージ交換サブルーチン	3-1
3-1 メッセージ交換ライブラリー	3-2
3-2 環境管理	3-4
3-3 集団通信	3-10
3-3-1 集団通信サブルーチンの一般的な規則	3-10
3-3-2 集団通信サブルーチンと同期	3-11
3-3-3 集団通信サブルーチンの種類	3-12
3-3-4 MPI_BCAST	3-13
3-3-5 MPI_GATHER	3-17
3-3-6 MPI_REDUCE	3-19
3-3-6-1 MPI_MAXLOC/MPI_MINLOC	3-22
3-3-6-2 ユーザーが定義する演算	3-23
3-4 1対1通信	3-28
3-4-1 1対1通信の全体の流れ	3-28
3-4-2 ブロッキング通信と非ブロッキング通信	3-29
3-4-3 1方向のみの通信	3-33
3-4-4 双方向通信	3-34
3-4-5 1対1通信に関する補足	3-38
3-5 派生データ型	3-39
3-5-1 派生データ型とは	3-39
3-5-2 派生データ型に関する主なMPIのサブルーチン	3-41
3-5-3 部分配列を表す派生データ型の作成	3-44
3-5-3-1 MPI_TYPE_CREATE_SUBARRAY	3-45
3-5-3-2 ユーティリティ・サブルーチン	3-46
3-5-4 構造体を表す派生データ型の作成	3-52
3-6 新グループの作成	3-53
3-7 C言語のMPI関数	3-54
3-8 MPI-2	3-56
3-8-1 MPI_IN_PLACE	3-56

## 第4章 プログラムの並列化の方法 .....4-1

4-1	パフォーマンスに関する考慮点	.....4-2
4-1-1	アムダールの法則	.....4-2
4-1-2	パフォーマンスを向上させるための考慮点	.....4-3
4-2	どんなところが並列化できるのか	.....4-7
4-3	計算部分の並列化パターン	.....4-10
4-4	並列化に伴う入力部分の修正	.....4-16
4-4-1	並列化に伴う入力部分の修正	.....4-16
4-4-2	並列化に伴う出力部分の修正	.....4-20
4-4-3	入出力に関するその他の考慮点	.....4-23
4-5	ループの分割方法	.....4-24
4-5-1	ループ反復と配列	.....4-24
4-5-2	配列の分割方法	.....4-25
4-5-3	ブロック分割でロードバランスが不均等になる例	.....4-27
4-5-4	ブロック分割	.....4-30
4-5-5	サイクリック分割	.....4-33
4-5-6	ブロック・サイクリック分割	.....4-35
4-5-7	配列の縮小	.....4-37
4-5-7-1	配列の縮小方法	.....4-38
4-5-7-2	配列の縮小に伴う変更(I/O以外の部分)	.....4-44
4-5-7-3	配列の縮小に伴う変更(I/O部分)	.....4-46
4-5-8	多重ループの分割	.....4-49
4-6	並列化とメッセージ交換	.....4-53
4-6-1	メッセージ交換はどんなときに必要か	.....4-53
4-6-2	1次元差分法	.....4-54
4-6-3	集団通信サブルーチンの代替	.....4-59
4-6-3-1	MP I _ G A T H E R V の代替	.....4-59
4-6-3-2	MP I _ A L L G A T H E R V の代替	.....4-62
4-6-3-3	MP I _ A L L T O A L L V の代替	.....4-65
4-6-4	合計(内積)/最大(最小)	.....4-68
4-6-5	特定の配列要素の参照	.....4-70
4-6-6	重ね合わせ	.....4-71
4-6-6-1	データの収集	.....4-71
4-6-6-2	間接アドレスを使用した計算	.....4-73
4-6-7	パイプライン法(ICCG法の前進消去/後退代入などで使用)	.....4-75
4-6-8	ツイスト分割法(ADI法などで使用)	.....4-79
4-6-9	一次逐次演算(漸化式などで使用)	.....4-84
4-6-10	受信するデータ量が不明な場合の通信	.....4-87
4-7	特殊な並列化方法	.....4-89
4-7-1	MPMDモデルでの並列化	.....4-89
4-7-1-1	実行方法	.....4-89
4-7-1-2	連成解析	.....4-90
4-7-1-3	マスター・スレーブ方式	.....4-91
4-7-2	フラットMPIとハイブリッド並列	.....4-94
4-8	並列化の手順と考慮点	.....4-96
4-8-1	並列化の手順	.....4-96
4-8-2	エラーかなとおもったら	.....4-103
4-8-3	パフォーマンスの測定と評価	.....4-106

第5章 並列化の応用例	5-1
5-1 差分法徹底攻略	5-2
5-1-1 2次元目でブロック分割した場合	5-3
5-1-2 1次元目でブロック分割した場合	5-4
5-1-3 1, 2次元目でブロック分割した場合	5-5
5-1-4 1, 2次元目でブロック分割した場合(斜めの要素を参照)	5-8
5-1-5 配列を縮小する場合	5-11
5-1-6 その他の考慮点	5-12
5-2 有限要素法	5-13
5-2-1 有限要素法の並列化	5-13
5-2-2 並列化の方法1	5-15
5-2-3 並列化の方法2	5-19
5-3 LU分解	5-27
5-4 ICCG法	5-33
5-4-1 パイプライン法による並列化	5-33
5-4-2 パラレル・ブロック・オーダリング法による並列化	5-46
5-5 マルチフロントアル法	5-58
5-6 SOR法	5-59
5-6-1 SOR法と並列性	5-59
5-6-2 プログラム例(ゼブラ法)	5-61
5-6-3 その他の解法	5-64
5-6-4 各解法の具体的な並列化方法	5-67
5-7 モンテカルロ法	5-73
5-8 個別要素法/分子動力学法	5-75
第6章 並列版数値計算ライブラリー	6-1
6-1 並列版数値計算ライブラリー	6-2
第7章 本当に並列化の必要があるか?	7-1
7-1 単体チューニング	7-2
7-2 並列処理と分散処理	7-3
7-3 並列化とワークロード	7-5
付録 MPIサブルーチン一覧	A-1
環境管理サブルーチン MPI_INIT	A-2
環境管理サブルーチン MPI_FINALIZE	A-3
コミュニケータに関するサブルーチン MPI_COMM_SIZE	A-4
コミュニケータに関するサブルーチン MPI_COMM_RANK	A-5
環境管理サブルーチン MPI_ABORT	A-6
環境管理フアункション MPI_WTIME	A-7
集団通信サブルーチン MPI_BCAST	A-8
集団通信サブルーチン MPI_SCATTER	A-10

集団通信サブルーチン	MPI_SCATTER	.....	A-12
集団通信サブルーチン	MPI_GATHER	.....	A-14
集団通信サブルーチン	MPI_GATHERV	.....	A-16
集団通信サブルーチン	MPI_ALLGATHER	.....	A-18
集団通信サブルーチン	MPI_ALLGATHERV	.....	A-20
集団通信サブルーチン	MPI_ALLTOALL	.....	A-22
集団通信サブルーチン	MPI_ALLTOALLV	.....	A-24
【MPI-2】集団通信サブルーチン	MPI_ALLTOALLW	.....	A-26
集団通信サブルーチン	MPI_REDUCE	.....	A-28
集団通信サブルーチン	MPI_ALLREDUCE	.....	A-30
集団通信サブルーチン	MPI_SCAN	.....	A-32
【MPI-2】集団通信サブルーチン	MPI_EXSCAN	.....	A-34
集団通信サブルーチン	MPI_REDUCE_SCATTER	.....	A-36
集団通信サブルーチン	MPI_OP_CREATE	.....	A-38
集団通信サブルーチン	MPI_BARRIER	.....	A-40
【MPI-2】MPI_IN_PLACEの使い方	MPI_SCATTER	.....	A-42
【MPI-2】MPI_IN_PLACEの使い方	MPI_SCATTERV	.....	A-43
【MPI-2】MPI_IN_PLACEの使い方	MPI_GATHER	.....	A-44
【MPI-2】MPI_IN_PLACEの使い方	MPI_GATHERV	.....	A-45
【MPI-2】MPI_IN_PLACEの使い方	MPI_ALLGATHER	.....	A-46
【MPI-2】MPI_IN_PLACEの使い方	MPI_ALLGATHERV	.....	A-47
【MPI-2】MPI_IN_PLACEの使い方	MPI_REDUCE	.....	A-48
【MPI-2】MPI_IN_PLACEの使い方	MPI_ALLREDUCE	.....	A-49
【MPI-2】MPI_IN_PLACEの使い方	MPI_SCAN	.....	A-50
【MPI-2】MPI_IN_PLACEの使い方	MPI_REDUCE_SCATTER	.....	A-51
1対1ブロッキング通信サブルーチン	MPI_SEND	.....	A-52
1対1ブロッキング通信サブルーチン	MPI_RECV	.....	A-54
1対1ブロッキング通信サブルーチン	MPI_SENDRECV	.....	A-56
1対1非ブロッキング通信サブルーチン	MPI_ISEND	.....	A-58
1対1非ブロッキング通信サブルーチン	MPI_IRECV	.....	A-60
1対1非ブロッキング通信サブルーチン	MPI_WAIT	.....	A-62
1対1通信サブルーチン	MPI_GET_COUNT	.....	A-64
派生データ型に関するサブルーチン	MPI_TYPE_SIZE	.....	A-66
派生データ型に関するサブルーチン	MPI_TYPE_CONTIGUOUS	.....	A-67
派生データ型に関するサブルーチン	MPI_TYPE_VECTOR/HVECTOR	.....	A-68
【MPI-2】派生データ型に関するサブルーチン			
派生データ型に関するサブルーチン	MPI_TYPE_CREATE_INDEXED_BLOCK	.....	A-70
派生データ型に関するサブルーチン	MPI_TYPE_INDEXED/HINDEXED	.....	A-72
【MPI-2】派生データ型に関するサブルーチン	MPI_TYPE_STRUCT	.....	A-74
派生データ型に関するサブルーチン	MPI_TYPE_CREATE_SUBARRAY	.....	A-76
【MPI-2】派生データ型に関するサブルーチン	MPI_TYPE_COMMIT	.....	A-78
派生データ型に関するサブルーチン	MPI_TYPE_CREATE_RESIZED	.....	A-79
通信データ型に関するサブルーチン	MPI_COMM_SPLIT	.....	A-80

(このページは空白です。)

# 第 1 章

## 並列プログラミングとは

本章では、分散メモリー型並列コンピュータの使用形態と、並列プログラミングの基本的な動作について説明します。

### ★注意★

- 本書の内容はときどき追加、更新しております。最新版は下記からダウンロード可能です。  
<http://accr.riken.jp/HPC/training/text.html>
- 本書内の「～はマシン環境に依存します」という記述は、その動作がマシンのハードウェア、OS、コンパイラなどの種類によって異なる可能性がある事を意味します。
- 本書のプログラム例はFortranで書かれています。C言語を使用していてFortranをご存じない方は以下の点に注意して下さい。
- 本書のプログラム例の中で特に宣言をしていない変数のデータ型は、先頭文字がI, J, K, L, M, Nで始まっている変数は整数、それ以外の変数は実数です。
- 実数が単精度(4バイト)か倍精度(8バイト)かは、プログラム内に「IMPLICIT REAL\*8(A-H,0-Z)」が指定されているかどうかによって以下のように決まります。

IMPLICIT REAL*8(A-H,0-Z)	なし	あり
REAL A	単精度	単精度
REAL*8 B	倍精度	倍精度
DIMENSION C(10)	単精度	倍精度
何も宣言されていない変数D	単精度	倍精度

- 本書のプログラムは原則として大文字で書かれていますが、Fortranでは大文字と小文字の区別がないので、どちらで書いても(あるいは混在しても)構いません。例えば変数AA, Aa, aA, aaは同一です。
- Fortranでは配列は1から始まります。  
【Fortran】 REAL A(2) : A(1), A(2)  
【C言語】 float a[2] : a[0], a[1]
- 多次元配列がメモリー上に配置される順番は、Fortranでは左の添字が先に動く順番になります。  
【Fortran】 REAL A(2,2) : A(1,1), A(2,1), A(1,2), A(2,2)  
【C言語】 float a[2][2] : a[0][0], a[0][1], a[1][0], a[1][1]

## 1-1-1 分散メモリー型並列計算機とは

### ■ 並列計算機の種類

#### (1) 単体計算機

単一CPUの計算機(以後単体計算機)は、図1-1-1(上段)に示すように、1つのCPUとメモリーから構成され、UNIXやLinuxなどのオペレーティングシステム(以下OS)が1つ存在してシステムを管理します。実行を開始した単体プログラムは単体プロセスと呼ばれ、メモリーにロードされます。図1-1-1(下段)に示すように、単体プログラムは、並列化して実行することが不可能な部分①、②(例えばI/O)と、可能な部分③、④(例えばD0ループの各反復で計算が独立している場合の1反復目と2反復目)で構成されています。

#### (2) 分散メモリー型並列計算機

並列計算機には、分散メモリー型と共有メモリー型があります。分散メモリー型は、図1-1-2(上段)に示すように、複数の単体計算機を、計算機間でデータの通信を行うためのネットワークで結合した計算機で、各CPUにそれぞれメモリーとOSが存在します。各CPUはお互いに完全に独立に実行しています。

分散メモリー用に並列化されたプログラムを2CPUで実行した場合の動作を図1-1-2(下段)に示します。

● ジョブが開始すると、並列化されていない①と②の部分(基本的に)全CPUで実行されます。

● 並列化されている③をプロセス0が、④をプロセス1が実行します。従ってこの部分の計算時間は理想的には半分になります。並列化された部分では、必要ならば、⇄に示すようにデータの通信を行います。

このように、並列ジョブの開始から終了までの間、プロセスの数はずっと2つです。つまり、並列化されている部分もされていない部分も、2つのCPUで並列に動作していることに注意して下さい。言いかええると、通信が可能な単体ジョブが、たまたま2つのCPUで同時に実行を開始したと考えることもできます。

#### (3) 共有メモリー型並列計算機

共有メモリー型では、図1-1-3(上段)に示すように、複数のCPUが1つのメモリーを共有します。またOSは1つだけ存在し、システム全体を管理します。なお、共有メモリー型並列計算機のことをSMP(Symmetric Multi-Processing)と呼ぶこともあります。

共有メモリー用に並列化されたプログラムを2CPUで実行した場合の動作を図1-1-3(下段)に示します。

● ジョブが開始すると、まず1つの並列プロセスが動作し、並列化されていない①を実行します。

● 並列化された部分に到達すると、並列プロセスはスレッドを1つ生成し、自分自身もスレッドとなります。スレッド0は③を、スレッド1は④を処理します。この部分の計算時間には半分になります。なお、スレッドというのにはUNIXの用語ですが、本書のレベルではプロセスと同じような物だと考えて下さい。

● ①と②の処理が終了すると、スレッド0は並列化されていない②を処理し、スレッド1はスレッド0が動的に変化します。このように共有メモリー型マシンでは、1つの並列ジョブの途中でスレッドの数が動的に変化します。

なお、分散メモリー型、共有メモリー型、共有メモリー型並列計算機で、単体プログラムを実行する事も勿論可能です。

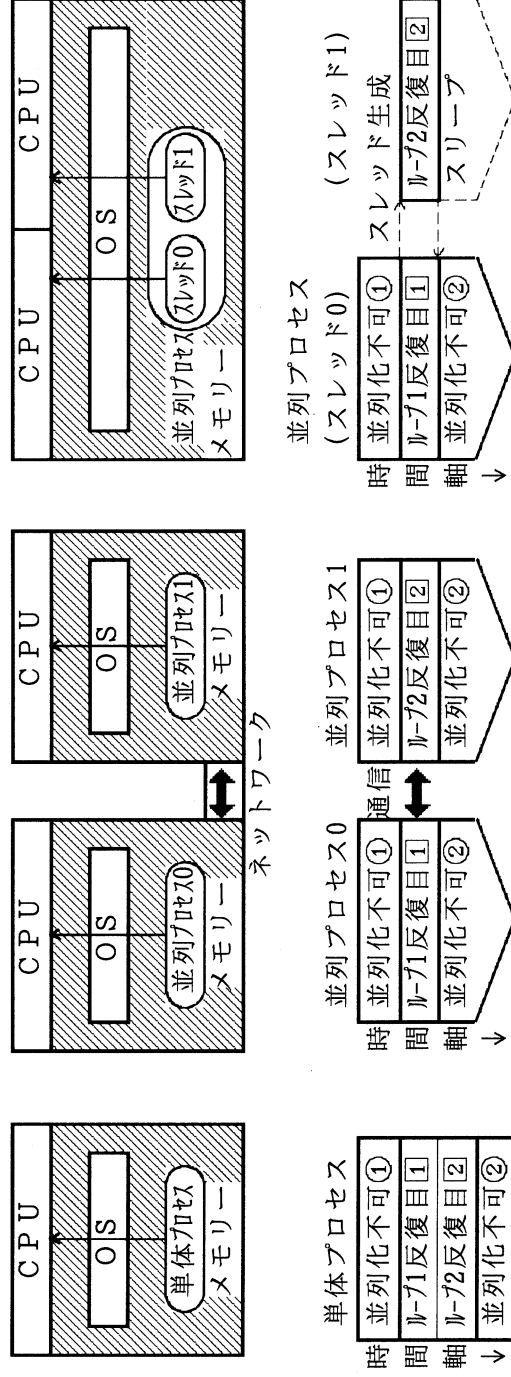


図1-1-1 単体計算機

図1-1-2 分散メモリー型並列計算機

図1-1-3 共有メモリー型並列計算機



■ 分散メモリ型並列計算機

以後、分散メモリ型並列計算機での並列化方法について説明します。分散メモリ型の並列計算機は、図1-1-4(1)に示すように、2台の単体の計算機をプログラム間でデータのやりとりができるようにネットワークで結合したのが最初の形態です。その後並列計算機が普及するにつれ、図1-1-4(2)に示すように、各CPUが並列計算機用の専用ラックに入ったタイプが出現しました。

一方図1-1-4(3)に示す共有メモリ型並列計算機では、通常は共有メモリ型の並列プログラムを実行しますが、マシン環境によっては(本書で説明する)分散メモリ型の並列プログラムを実行することもできます(4-7-2節参照)。

最近では、図1-1-4(4)に示すように、複数の共有メモリ型並列計算機をネットワークで結合した分散メモリ型並列計算機が増えてきています。

並列計算機では、各CPUをノード、あるいはPE(Processing Element)と呼ぶこともあります(本書では主に「ノード」を使用します)。並列プログラムでは、ノード間でデータを通信する速度が遅いと並列化の効果が出ないので、通常高速なネットワークを使用します。

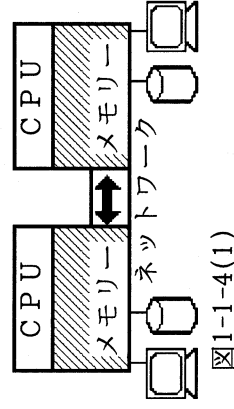


図1-1-4(1)

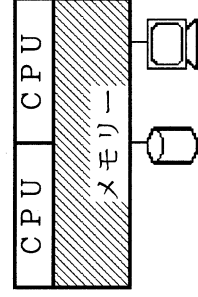


図1-1-4(3)

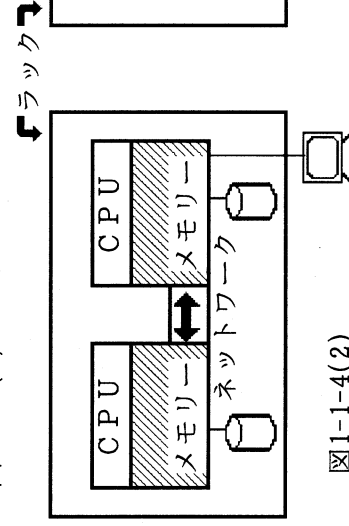


図1-1-4(2)

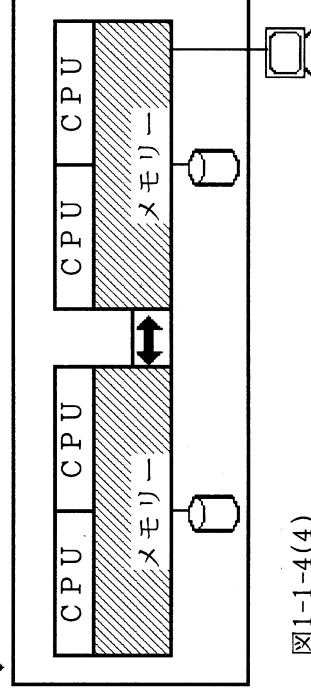


図1-1-4(4)

■ 並列計算機の使用形態

並列計算機には、以下のような使用形態があります。

- (1) 並列ジョブを実行する。
- (2) 単体ジョブしか実行しない。すなわち単体計算機の集合体として使用する。
- (3) 両者を混合させて使用する。

並列ジョブを実行する場合、既に並列化されている市販のパッケージを使用する場合と、自分でプログラムを並列化する場合があります。以下に並列化されている市販パッケージの例を示しますが、使用できるかどうかはお使いのマシンに依存します。

- 構造解析 : ABAQUS, LS-DYNA3D, PAM-CRASH, MARC, NASTRAN, RADIOSS
- 流体 : STAR-CD, FLUENT, STREAM, SCRYU
- 計算化学 : DISCOVER, HONDO, GAMESS, GAUSSIAN

一方自分でプログラムを並列化する場合、大きく次の3つの目的があります。

- 経過時間を短縮するため
- 計算に大容量のメモリが必要なため
- 並列アルゴリズムを研究するため

## 1-2 並列プログラミング・モデル

### ■ SPMDモデルとMPMDモデル

図1-2-1(1)では、単体プログラムをコンパイル・リンクしてロードモジュールa.outを作成し、3つのノードで同時に実行を開始しています。各単体プロセスの自身は同じですが、当然ながらお互いに独立に動作します。また各プロセスは同じ変数Aを使用していますが、異なるプロセスなので全く関係はありません。

次に並列プログラムの説明をします。並列プログラミング・モデルとして、SPMD(Single Program

Multiple Data)モデルとMPMD(Multiple Program Multiple Data)モデルの2種類があります。

SPMDモデルでは、図1-2-1(2)に示すように、1つの並列プログラムから作成された同じロードモジュールa.outが各ノードにロードされ、並列プロセスとして同時に実行を開始します。各プロセスにはランクという識別子(ID)(3プロセスの場合、**0,1,2**)が付きます。各プロセスの自身は同じですが、お互いに独立に動作します。また各プロセスは同じ変数Aを使用していますが、異なるプロセスなので全く関係はありません。つまり、図1-2-1(1)と同じ状態であると考えると分かります。図1-2-1(1)と唯一異なるのは、必要であれば各プロセス間で通信を行ってデータを交換することができる点です。なお、SPMDモデルでは各プロセスの自身は同じですが、全く同じ動作をするわけではなく、一般に異なる動作をします。一方MPMDモデルでは、図1-2-1(3)に示すように、構造解析、流体解析、流体解析という全く異なる複数の並列プログラム(ロジック、変数名も全く関係がない)から作成された異なるロードモジュールa.outとb.outが各ノードにロードされ、並列プロセスとして同時に実行を開始します。これは前述のように、構造解析、流体解析という単体プログラムを、本例では2本ずつ計4本同時に実行したのと同じ状態であると考えて下さい。唯一異なるのは、必要であれば各プロセス間で通信を行ってデータを交換することができる点です。本例のような処理を連続解析と呼び、原子力の解析や、地球規模の解析(大気+海洋)などで用いられます。

図1-2-1(4)はマスター・スレーブ(またはマスター・ワーカー)と呼ばれる方式で、1つのノードで親プロセスa.outを、その他全てのノードで同一の子プロセスb.outを実行します。親プロセス(係長)は子プロセス(平社員)に仕事を与え、子プロセスはその仕事を終了したら親プロセスに報告し、親プロセスは次の仕事の子プロセスに与えます。本例はマスターとスレーブを別プログラムで作成しているので、MPMDモデルです。

ほとんどのプログラムはSPMDモデルで並列化するので、以後はSPMDモデルでの並列化方法を説明し、MPMDモデルは4-7-1節で説明します(マシン環境によっては、MPMDモデルを実行できない場合もあります)。

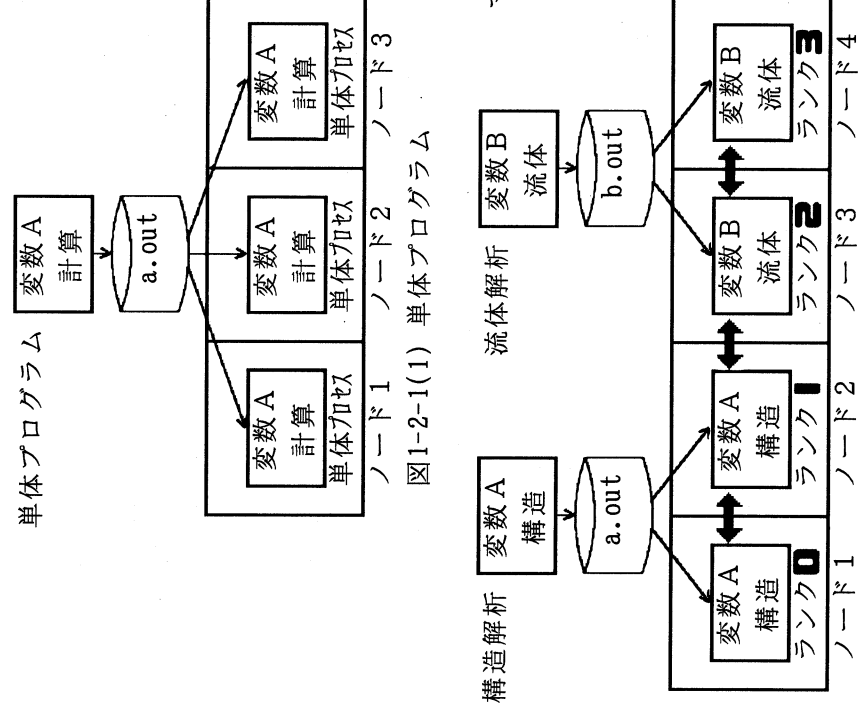


図1-2-1(1) 単体プログラム

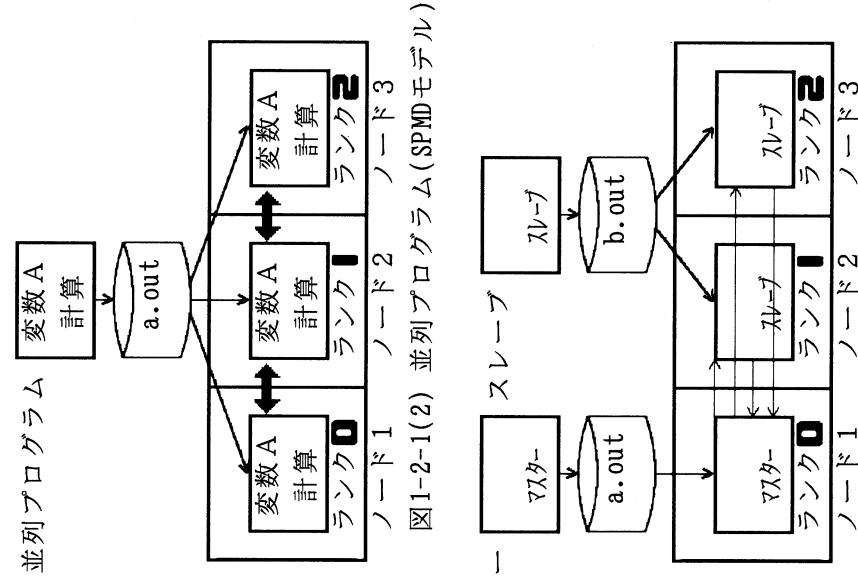


図1-2-1(2) 並列プログラム (SPMDモデル)

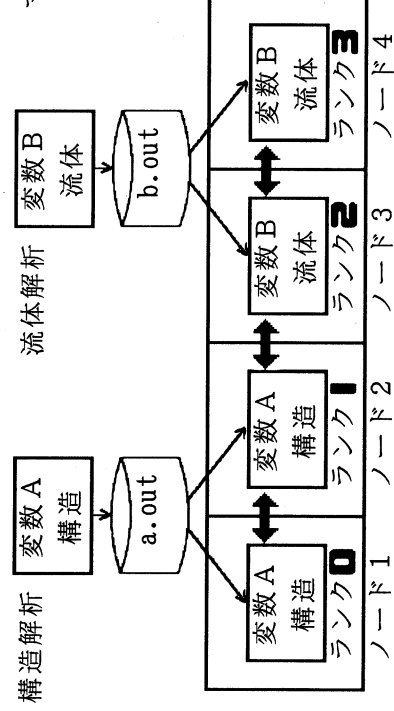


図1-2-1(3) 並列プログラム (MPMDモデル)

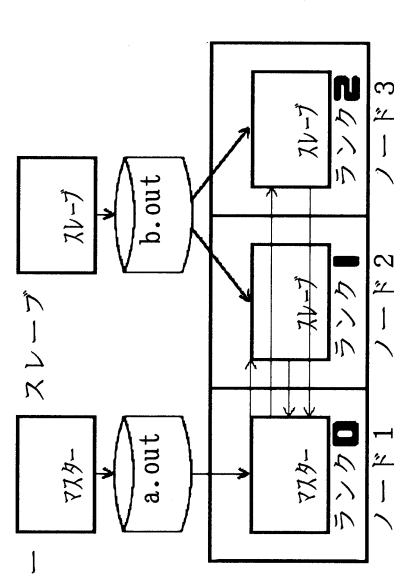


図1-2-1(4) 並列プログラム (MPMDモデル)

## ■ SPMDモデルの動作例(1)

SPMDモデルでは、1つの並列プログラムからロードモジュールを作成し、全ノードで同じ並列プロセスを実行しますが、各プロセスは全く同じ動作をするわけではなく、一般に異なる動作をします。つまり、1つの並列プログラム内に、全プロセスの動作が記述されています。これを簡単な例で説明します。

## (1) 単体での動作

一般のプログラムはマクロに見ると、『データの入力』、『処理』、『結果の書き出し』の3つの部分に分かれています。これを単純化したプログラムを単体で実行した例を図1-2-2(1)に示します。

- ①で入力データを読み込み、これを大きさ6の配列に入れます。
- ②で、配列内で処理するデータの最初と最後の位置をISとIEに設定します。
- ③でISからIEまでの各データ(つまり6つのデータ全て)に対して処理を行います。処理前と処理後のデータを、○△□と●▲■で表します。③はD0文(Fortranの場合)、for文(C言語の場合)に相当します。
- 最後に④で、処理したデータをファイルに書き出します。

## (2) SPMDモデルで並列化した場合の動作

図1-2-2(1)のプログラムをSPMDモデルで並列化し、3つのプロセスをノード2, 4, 6の3つのノードで実行する例を図1-2-2(2)に示します。

- 並列化したプログラムのロードモジュールがa.outであるとし、並列ジョブを開始すると、a.outが各ノードにロードされ、並列プロセスとして動き始めます。例えば3つのプロセスを実行した場合、前述のように、各プロセスにはランクと呼ばれる識別子(ID)(本例では**0, 1, 2**)が付けられます。このランクの値が各プロセスで異なるため、SPMDモデルで、同じ自身のプロセスに異なる動作をさせることができます。
- 図中でプログラムが3つ示してありますが、SPMDモデルなので中身は同一であることに注意して下さい。①で全プロセスがそれぞれ同じ入力データを大きさ6の配列に読み込みます。分散メモリー型マシンなので、データを読み込む配列は各ノードのローカルメモリー上にあり、お互いに全く関係がないのは言うまでもありません。なお、入力データを読み込む方法として、1つのプロセスのみが代表して読み込み、それを他のプロセスに送信する方法もあります。
- ②で、自分のランクを取得します。例えば自分がランク**0**のプロセスであれば**0**が戻ります。なお、①と②の順序は逆でもかまいません。
- ③で、処理するデータの最初と最後の位置をISとIEに設定します。ランクによって異なる部分(枠で囲んだ部分)が実行されるため、異なるISとIEの値が設定されます。
- ④で処理を行います。③で自分が担当するISとIEの値が設定されているので、各プロセスはそれぞれ自分の担当部分のみを処理します。例えばランク**0**のプロセスは●の部分のみを処理します。
- ④が終了した後、各プロセスは自分が処理した部分のデータを⑤で(例えば)ランク**0**のプロセスに送信します。
- ⑥で、ランク**0**のプロセスが代表して結果の書き出しを行います。

実際のプログラムをSPMDモデルで並列化した場合も、マクロに見るとこの例のようになります。ただし、実際のプログラムでは通常④の部分でデータの通信が発生します。

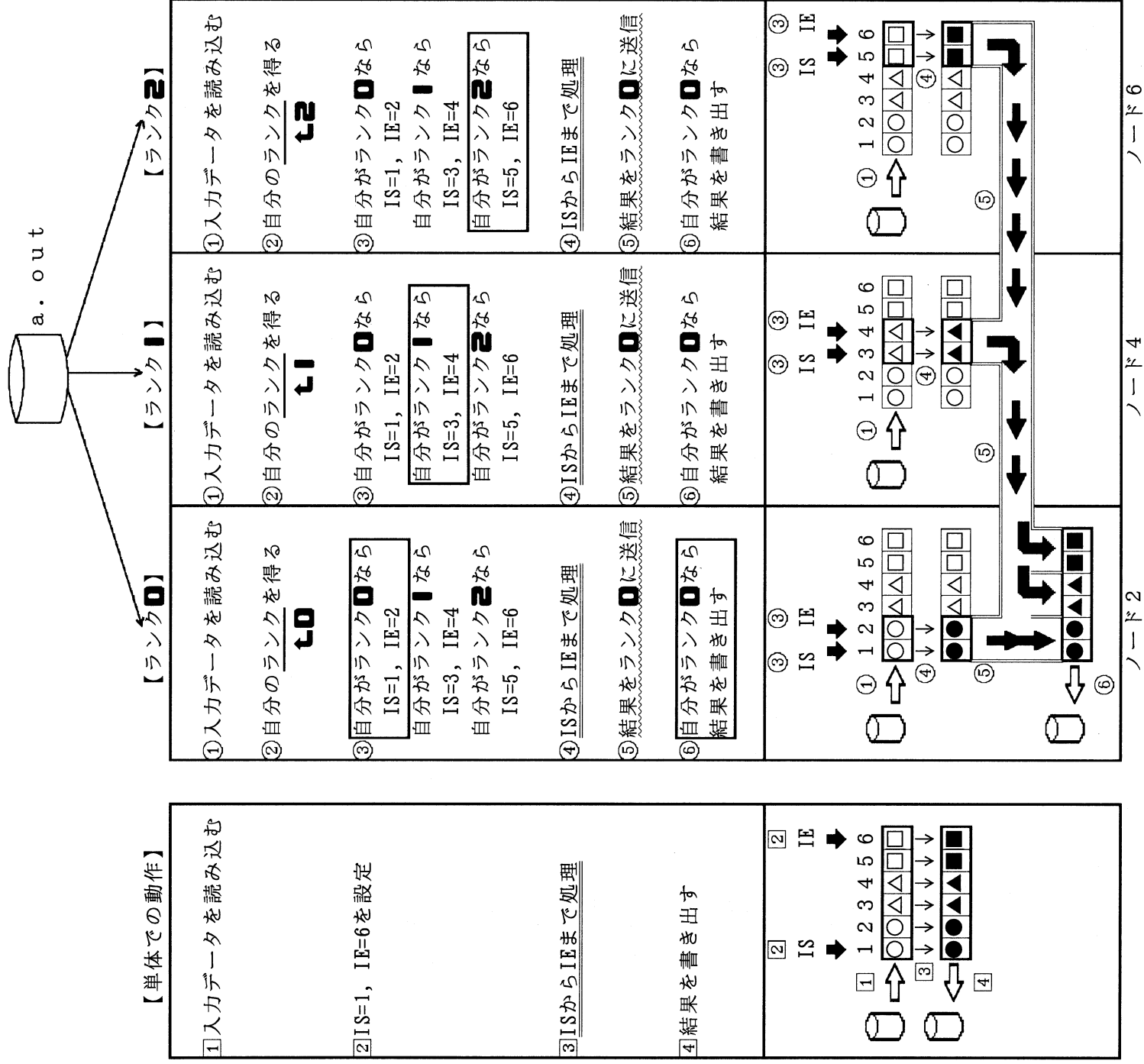


図1-2-2(1) 単体での動作

図1-2-2(2) SPMDモデルで並列化した場合の動作

## ■ 並列化とメッセージ交換

### (1) 並列化

3章以降で、ノード間でデータの通信(図1-2-2(2)の⑤)を行うメッセージ交換サブルーチンの使用方法について説明していきます。メッセージ交換サブルーチンの使用方法を学ぶにつれ、あたかも『メッセージ交換 = 並列化』であると考えようになりがちですが、これは誤りです。このことを説明する前に、まず本当の意味の並列化について定義します。

図1-2-2(2)から分かるように、各プロセスのプログラムは最初から最後まで並列に動作します。しかし④以外の部分は、単体で動作した場合と計算量が変わらない(または並列化に伴って新規に追加された)部分なので、本当の意味の並列化とは言えません。本当の意味の並列化と言えるのは、計算量が1/3に減少した④の部分のみです。

そこで本書では、本当の意味の並列化を次のように定義します。

本当の意味の並列化とは、計算量を $1/n$ に縮小すること(ノード数が $n$ 台の場合)。

### (2) メッセージ交換

本当の意味の並列化とメッセージ交換(図1-2-2(2)の⑤)はどのような関係になるのでしょうか。例えば図1-2-2(2)では、本当の意味の並列化を行った④が終了した直後、各プロセスは自分が処理した1/3のデータのみを持っており、他プロセスが処理した部分のデータは、(自分とは別のプロセスなので)直接参照/更新することはできません。

このことによつて、本当の意味の並列化を行った部分以後の処理に矛盾(副作用)をきたさないの**であれば**、メッセージ交換を行う必要はありません。しかしこの例では、図1-2-2(1)と同じ結果を得るために、各プロセスが処理したデータを、⑤で(例えば)ランク0のプロセスに集める必要がある**ため**、このときメッセージ交換を行います

まとめると、本当の意味の並列化とメッセージ交換の関係は以下ようになります。

本当の意味の並列化に伴う矛盾(副作用)を解消するために、必要最小限、仕方なしに行うのがメッセージ交換。

## ■ SPMDモデルの動作例(2)

SPMDモデルによる並列化のイメージをつかむために、動作例をもう一つ示します。まず単体での動作を図1-2-3(1)に示します。図中の大きさ9の配列に1～9の値が入っており、それを総合計した結果をディスプレイに表示します。

これをSPMDモデルで並列化した例を図1-2-3(2)に示します。

- 図1-2-2(2)と同様に、自分のランクから、自分の担当するデータの最初と最後の位置を決め、これをIS、IEとします。
- 各プロセスは自分が担当する部分の合計(小計)を求めます。この部分は計算量が1/3に減少したので、『本当の意味の並列化』になります。
- この時点で、各プロセスが持っている値は自分が計算した部分の小計です。これが並列化に伴う矛盾(副作用)です。もちろんこのままであればそれで終わりですが、この例では小計の合計、すなわち総合計を求めたいので、メッセージ交換を行います。これによって、並列化に伴う矛盾(副作用)を解消することが出来ます。
- この場合、図1-2-2(2)と同様に、各プロセスが求めた小計(6, 15, 24)を一度ランク0のプロセスに集め、その後ランク0のプロセスが代表して(6+15+24=45)を計算するという方法もあります。しかしこれでは面倒なので、図1-2-3(2)では、『通信しながら演算を行う』タイプのメッセージ交換サブルーチンを使用しました。これを使用すると、小計(6, 15, 24)に対して通信しながら演算が行われ、結果が例えばランク0のプロセスに入ります。
- 最後に、ランク0のプロセスが代表して総合計をディスプレイに表示します。

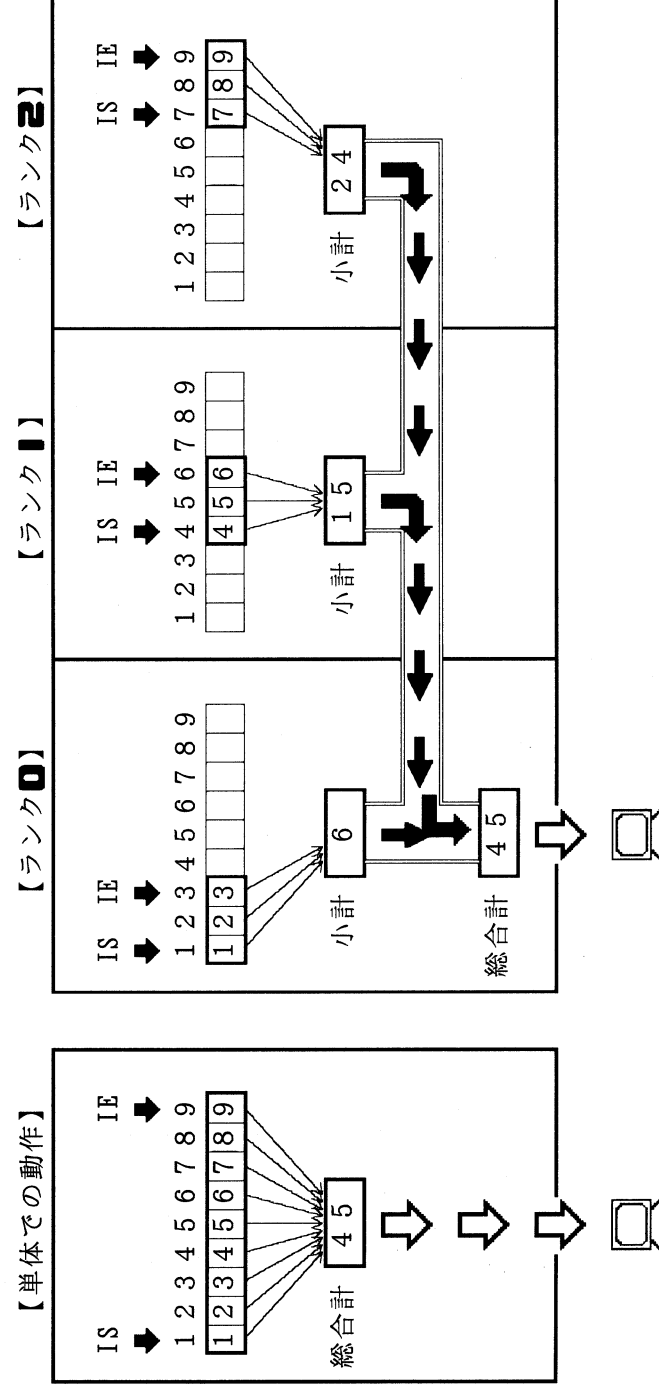


図1-2-3(1) 単体での動作

図1-2-3(2) SPMDモデルで並列化した場合の動作

## 第2章 並列ジョブの実行方法

本章では、簡単なプログラムを用いて、並列ジョブのコンパイル・リンク・実行方法について説明します。

## 2-1 並列ジョブの実行方法

### ■ コンパイルから実行までの流れ

簡単なプログラムを用いて、並列ジョブのコンパイルから実行までの一連の流れを見てみましょう。なお、具体的な操作方法はマシン環境によって異なりますので、お使いのマシンのユーザーズガイド等を参照して下さい。

- 図2-1-1の①に示す簡単なプログラムhello.fを並列に実行してみます。なお④のプログラムは(後述する)MPIのサブルーチンMPI\_INITやMPI\_FINALIZEなどが指定されていないため、マシン環境によっては正常に動作しなかったり異常終了することがあります。
- ②でhello.fをコンパイル・リンクすると、③でロードモジュールa.outが作成されます。ここでは、a.outは各ノードからアクセスできる共有ディスクに作成されます。なお、コンパイル・リンクの方法はマシン環境によって異なります。
- ④で、並列ジョブa.outを2プロセスで実行します。すると⑤でa.outが各CPUにロードされ、ランク0、ランク1のプロセスとして実行を開始します。なお、並列ジョブの実行方法および、どのプロセスがどのCPUにロードされるかは、マシン環境によって異なります。
- ⑥と⑦のPRINT文が並列に実行され、⑧と⑨が画面に表示されます。⑧と⑨の表示順序は早い者順になります。なお、この例のように複数のプロセスが同時に標準出力(WRITE(6)またはPRINT文)に書き出しを行った場合、出力行数が多かったり各行の文字数が多いと、各行が混ざって表示が乱れる事がありますので、実際のプログラムでは標準出力への書き出しは(例えば)ランク0のプロセスのみが行うようにします(4-4-2節参照)。

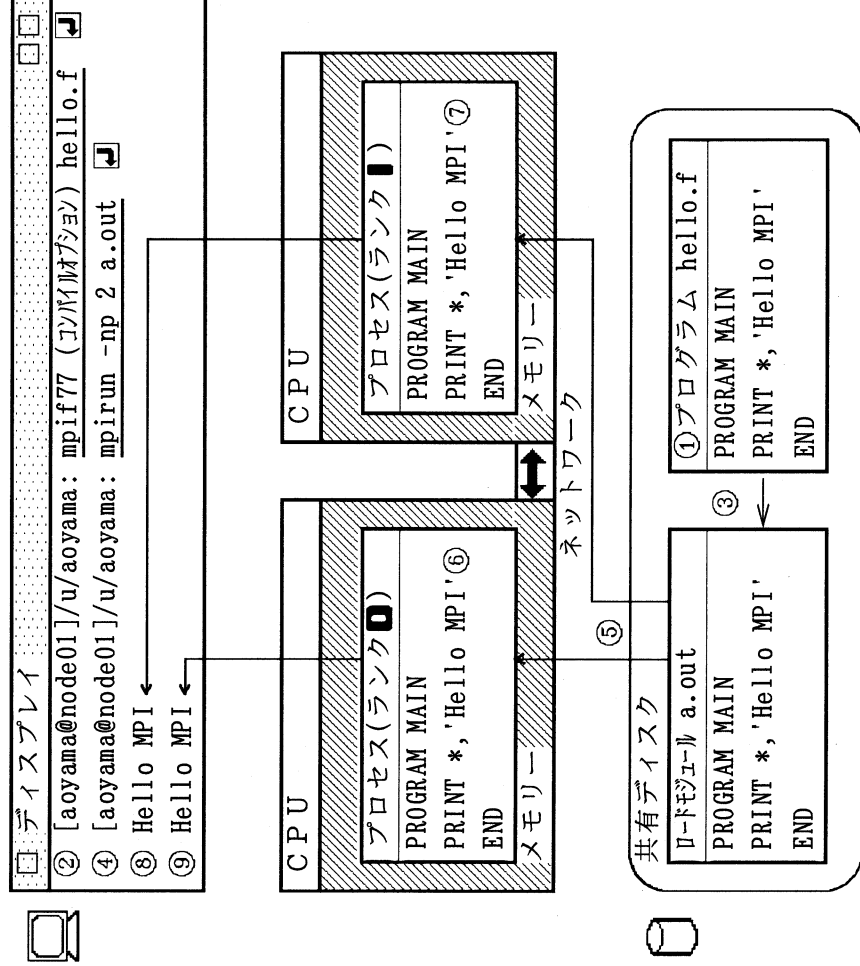


図2-1-1



■ 並列ジョブをサポートするソフトウェア

並列ジョブの開発や実行をサポートするために、以下のようなソフトウェアが提供されていることがありますが、マシン環境によって状況が異なるので、ここでは簡単に説明します。なお、私自身は(3)~(5)のツールは使ったことがありません。

(1) 並列ジョブ管理

並列ジョブを実行したとき、各プロセスをどのノードに割り当てるかなど、並列ジョブの実行に関する管理を行います。

(2) メッセージ交換ライブラリー

並列プログラムが、ノード間でデータの通信を行うためにコールする、通信ルーチンのライブラリーです。MPI(後述)というライブラリーが一般的に使われます。

(3) 並列デバッガ

並列プログラム用のdbxデバッガです。なお、私は並列プログラムのデバッグはFortranのPRINT文のみで行っています。

(4) 並列プロファイラー

プログラムに含まれる各サブルーチンのCPU時間の割合を調べるためのUNIXコマンドprof、gprofの並列版です。

(5) 可視化ツール

並列ジョブの活動状況を可視化して表示します。

(このページは空白です。)

## 第3章

### メッセージ交換サブルーチン

第1章ではメッセージ交換について、『本当の意味の並列化に伴う矛盾(副作用)を解消するために、必要最小限、仕方なしに行うのがメッセージ交換』などと、すげなく扱ってしまいました。メッセージ交換は英話で言えば単語に相当し、プログラムの並列化は英会話に相当します。従って、プログラムの並列化を行うためにはメッセージ交換サブルーチンの使い方をしっかり学ぶ必要があります。

## 3-1 メッセージ交換ライブラリー

(注) 本書中で、「マシン環境によって動作が異なります。」等と記載した部分は、並列計算機のハードウェア、ソフトウェア、コンパイラ、MPIライブラリー等によって動作の異なる(可能性のある)部分です。ただし、私は2種類の計算機しか使用した事がないので、上記の記載をした以外の部分でもマシン環境によって動作の異なることがあるかもしれません。

### ■ MPI概要

分散メモリー型のマシンで並列計算を行う際、メッセージ交換(通信)が必要になります。MPIが作成される以前に、PVM(Parallel Virtual Machine)などいくつかのメッセージ交換ライブラリーが存在しました。

しかし並列プログラムの互換性の観点から、ライブラリーの標準化の機運が高まり、米国のコンピュータメーカーや大学の研究者がMPIフォーラムという組織を結成し、既存のライブラリーの長所を集めたMPI(Message-Passing Interface)と呼ばれる基本仕様を1995年に作成しました(参考文献[15][27]参照)。この仕様をMPI-1と呼んでいます。

その後、MPI-1の拡張機能仕様であるMPI-2が作成されました。本書ではMPI-1の範囲で説明を行います。

MPI-1	MPI-2
基本仕様	拡張機能仕様

MPIはFortranおよびCからコールすることができます(本書では全てFortranで説明を行います)。

MPI-1に含まれているサブルーチンは図3-1-1に示すように6つに分類されます。そしてサブルーチンの総数は何と128個(処理系によって多少異なります)にもなります。では、この128個のサブルーチンを全てマスターしないと並列プログラムは作成できないのでしょうか？

サブルーチン分類	数
1対1通信	39
集団通信	16
派生データ型	13
コミュニケーション	30
プロセス・トポロジー	16
環境管理	14
合計	128

図3-1-1

答えは『No!』です。実はこれらの128個の中には、何のために使用するかよく分からないルーチンや、単に照会をするだけのサブルーチンも多く含まれており、コンピュータメーカーでMPIの処理系を開発する人や並列処理の研究者以外は、そのようなルーチンは知らなくてもよいと言っても過言ではありません。そして、よほど特殊な並列化を行わない限り、20個程度を知っていれば十分です(その中でも特によく使用するのは10個程度です)。そこで本書では、重要と思われる約20個のみを取り上げ、それ以外のサブルーチンは全て割愛することにしました。

■ 付録について

MPIで提供されている主なサブルーチンの使用方法を、巻末の付録にまとめました。また付録内の「機能」の「関連する節」の所に、各サブルーチンが出てくる本書内の節を示しました。

- 言うまでもありませんが、本書は正式マニュアルではありませんので、正式な使用方法は参考文献[15][27]を参照して下さい。
- 付録に掲載したのはFortran版のMPIルーチンの使用方法です。例えばFortran版で引数が整数になっているても、C言語版では別のデータ型で宣言しなければならぬ場合がありますので、C言語のユーザーの方は、参考文献[15][27]内のC言語の説明箇所を参照して下さい。
- Fortran版のMPIルーチンで引数が整数の場合、下記の下線部を付けて宣言しなければならぬ場合があります(例えば付録のMPI\_TYPE\_CREATE\_RESIZED参照)。

```
INTEGER(KIND=MPI_ADDRESS_KIND) N
CALL MPI_xxx(~,N,~)
```

■ 多次元配列のメモリー上での並び

多次元配列の一部のデータを通信するような場合、配列の要素がメモリー上でどのように並んでいるかを知っておく必要があります。

まず、本テキストでは、2次元配列の図は、図3-1-2(1)に示すように、1次元目を縦方向、2次元目を横方向に描いています。この2次元配列の各要素が、メモリー上でどのように並んでいるかを説明します。

Fortranの場合、2次元配列の各要素は、左側の添字が先に動く順番にメモリー上に並びます。例えば2次元配列A(3,4)は、メモリー上では図3-1-2(2)に示す順に並びます。図3-1-2(1)では、**■**、**▲**、**◆**、**●**の各要素は一見連続しているように見えますが、実際には図3-1-2(2)に示すように、メモリー上ではとびとびになっているため、それを意識して通信を行う必要があります。

3次元以上の配列の場合も、最も左側の添字が一番先に動く順番にメモリー上に並びます

一方C言語の場合、Fortranと仕様が逆になっており、2次元配列の各要素は、右側の添字が先に動く順番にメモリー上に並びます。例えば2次元配列a[3][4]は、メモリー上では図3-1-3(2)に示す順に並びます。

<b>■</b> A(1,1)	<b>▲</b> A(1,2)	<b>◆</b> A(1,3)	<b>●</b> A(1,4)
A(2,1)	A(2,2)	A(2,3)	A(2,4)
A(3,1)	A(3,2)	A(3,3)	A(3,4)

図3-1-2(1) (Fortran) 2次元配列A(3,4)

:
<b>■</b> A(1,1)
A(2,1)
A(3,1)
<b>▲</b> A(1,2)
A(2,2)
A(3,2)
<b>◆</b> A(1,3)
A(2,3)
A(3,3)
<b>●</b> A(1,4)
A(2,4)
A(3,4)
:

図3-1-2(2) (Fortran) A(3,4)のメモリー上の配置

<b>■</b> a[0][0]	<b>▲</b> a[0][1]	<b>◆</b> a[0][2]	<b>●</b> a[0][3]
a[1][0]	a[1][1]	a[1][2]	a[1][3]
a[2][0]	a[2][1]	a[2][2]	a[2][3]

図3-1-3(1) (C言語) 2次元配列a[3][4]

:
<b>■</b> a[0][0]
<b>▲</b> a[0][1]
<b>◆</b> a[0][2]
<b>●</b> a[0][3]
a[1][0]
a[1][1]
a[1][2]
a[1][3]
a[2][0]
a[2][1]
a[2][2]
a[2][3]
:

図3-1-3(2) (C言語) a[3][4]のメモリー上の配置

## 3-2 環境管理

本節では環境管理に関するサブルーチンについて説明します。なお、本書のFortranプログラム例は原則として大文字で書きますが、Fortranは大文字と小文字の区別がないので、小文字で書いても構いません。

まず図3-2-1(1)のプログラムを作成し、3プロセスで実行した結果を図3-2-1(2)に示します。図中でプログラムは3つあるように見えますが、SPMDモデルのためプログラムは1つです。

図3-2-1(1)のプログラムは、MPIを使用した最も簡単なプログラムであると言えます。以後各ステートメントについて説明します。

ランク0

```
PROGRAM MAIN
INCLUDE 'mpif.h'
  .. ①
(②の前にユーザープログラムが存在しても構いません)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR) .. ③
  ↳ 3
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR) .. ④
  ↳ 0
PRINT *, 'NPROCS = ', NPROCS, 'MYRANK = ', MYRANK .. ⑤
CALL MPI_FINALIZE(IERR)
  .. ⑥
(⑥の後にユーザープログラムが存在しても構いません)
END
```

ランク1

```
.. ①
しても構いません)
NPROCS, IERR) .. ③
  ↳ 3
MYRANK, IERR) .. ④
  ↳ 1
= ', MYRANK .. ⑤
  .. ⑥
しても構いません)
END
```

ランク2

```
.. ①
しても構いません)
NPROCS, IERR) .. ③
  ↳ 3
MYRANK, IERR) .. ④
  ↳ 2
= ', MYRANK .. ⑤
  .. ⑥
しても構いません)
END
```

図3-2-1(1)

```
[aoyama@node011]~/aoyama: mpirun -np 3 a.out
NPROCS = 3 MYRANK = 0
NPROCS = 3 MYRANK = 1
NPROCS = 3 MYRANK = 2
```

図3-2-1(2)

### ■ インクルードファイルmpif.h

図3-2-1(2)①のインクルードファイルmpif.hは「/usr」などのディレクトリの下に置かれており、置かれる場所や中身はMPIの処理系によって異なります。ユーザーはmpif.hの中身について知る必要はありませんが、参考までに例を図3-2-2(1)に示します。図中の「MPI\_xxx」という整数(例えば下線の「MPI\_COMM\_WORLD」)はMPIが内部的に使用する定数で、各定数の具体的な値がPARAMETER文で設定されています。ユーザーは「MPI\_xxx」の各定数を、MPIのサブルーチンの引数で指定しますが、定数名だけを知っていれば良く、PARAMETER文で設定されている具体的な値については関知する必要はありません。

なお、「MPI\_」、「PMPI\_」で始まる変数名、関数名、サブルーチン名をユーザープログラムで使用すると、MPIが使用する変数と偶然一致して誤動作する危険性がありますので、使用しないで下さい。

```
:
integer MPI_COMM_WORLD
parameter (MPI_COMM_WORLD=0)
integer MPI_INTEGER, MPI_REAL4, MPI_REAL, MPI_REAL8
parameter (MPI_INTEGER=18, MPI_REAL4=19, MPI_REAL=19, MPI_REAL8=20)
:
```

図3-2-2(1) mpif.hの中身の例

次にmpif.hの指定方法について説明します。図3-2-1(1)で、「CALL MPI\_～」となっていてサブルーチン(②,③,④,⑥)はMPIのサブルーチンです。プログラム内のメインルーチンとサブルーチンの内、MPIのサブルーチンを使用した全てのルーチンで、必ず①の「INCLUDE 'mpif.h'」を指定して下さい(必須)。なお、Fortranは大文字と小文字の区別がありませんが、「INCLUDE 'MPIF.H'」のように「mpif.h」の部分を大文字で書くこととエラーになるマシン環境もありますので、小文字で書いた方がよいでしょう。

指定する場所は、図3-2-2(2)の③に示すように、宣言文の先頭(ただしUSE文とIMPLICIT文の後)に置くのがよいでしょう。理由ですが、例えば⑥の下線部のMPI\_STATUS\_SIZE(意味は3-4節で説明します)のように、MPIが使用する定数をプログラム内で使用することがあり、MPI\_STATUS\_SIZEの値は前述のように④のファイル内のPARAMETER文で定義されているため、⑥は④の後に置く必要があるからです。

別の方法として、図3-2-2(3)の③のように、「INCLUDE 'mpif.h'」や並列化に関係するその他の変数をMODULE文(Fortran90の機能)の中に含める方法もあり、私はこの方法を用いています(4-8-1節参照)。

```
PROGRAM MAIN(またはサブルーチン)
USE xxxxx
IMPLICIT REAL*8(A-H,0-Z)
INCLUDE 'mpif.h'
通常の宣言文
INTEGER ISTATUS(MPI_STATUS_SIZE)
:
CALL MPI_xxx
:
```

図3-2-2(2)

```
MODULE PARA
INCLUDE 'mpif.h'
:
END
PROGRAM MAIN(またはサブルーチン)
USE PARA
IMPLICIT REAL*8(A-H,0-Z)
通常の宣言文
:
```

図3-2-2(3)

■ MPI\_INIT

まず付録のMPI\_INITの説明を一読して下さい....。

図3-2-1(1)の②はMPI環境の初期化処理を行うサブルーチンで、全てのMPIルーチンの一番最初に1回だけ、必ずコールする必要があります(必須)。MPI\_INITをコールする前にユーザープログラムが存在しても構いません。

MPI\_INITで指定する整数の引数 ierrror(図3-2-1(1)ではIERR)は、実は本サブルーチンに限らずMPIの全てのサブルーチンの最後の引数でも同様に指定します。引数IERRは必ず指定する必要がありますが、使用することはほとんどありません。IERRには、コールしたMPIのサブルーチンが正常に終了すれば「MPI\_SUCCESS」(具体的な値はインクルードファイルmpif.hの中で与えられています)が、エラーが発生した場合はそれ以外の値が戻ってきます。しかし、並列ジョブを実行していてMPIのサブルーチンでエラーが発生した場合、自動的にジョブは異常終了し、エラーメッセージが標準エラー出力に表示されます。従って、通常、MPIのサブルーチンをコールした後にユーザープログラムで引数IERRの値をいちいちチェックする必要はありません。しかし、必要がないからと言ってIERR自体を引数として指定しないと、「segmentation fault」等エラーになりますので、IERRは「オマジナイ」として必ず指定するようにして下さい。

なお、コールしたサブルーチンでエラーが発生した場合、IERRに戻る「MPI\_SUCCESS」以外の値と実際のエラーメッセージの対応は、「CALL MPI\_ERROR\_STRING」というサブルーチンをコールすることによって調べることができますが、本書では説明は省略します。

■ MPI\_FINALIZE

まず付録のMPI\_FINALIZEの説明を一読して下さい....。

図3-2-1(1)の⑥はMPI環境の終了処理を行うサブルーチンで、全てのMPIルーチンの一番最後に1回だけ、必ずコールする必要があります(必須)。MPI\_FINALIZEの後にユーザープログラムが存在しても構いません。

前述のMPI\_INITは、コールし忘れると通常実行時にエラーで異常終了しますが、MPI\_FINALIZEをコールせずに終了した場合、マシン環境によっては誤動作する(実行中の他のプロセスを強制的に終了させてしまう、あるいは実行するたびに計算結果が異なる再現性のないエラーとなる)ことがあります。そしてその原因が、よもやMPI\_FINALIZEのコールし忘れだとは、通常気付きません。従って、プログラムが終了する全ての箇所(STOP文またはEND文)で、必ず1度だけMPI\_FINALIZEをコールしてから終了するようにして下さい。

## ■ MPI\_COMM\_SIZE / MPI\_COMM\_RANK

MPI\_COMM\_SIZEとMPI\_COMM\_RANKは、図3-1-1の分類では「環境管理」でなく、「コミュニケーション」に属しますが、便宜的にここで説明します。

例えば「mpirun -np 3 a.out **□**」で並列ジョブを実行すると、図3-2-3に示すように3つのプロセスが作成され、各プロセスにはランク**0**~**2**のいずれかの値が付けられます。その際、この3つのプロセス(全プロセス)からなるグループが自動的に作成され、MPI\_COMM\_WORLDというグループ名(コミュニケーションと呼べれます)が付けられます。(正確には、MPI\_COMM\_WORLDはMPIが内部的に使用する定数(整数)で、前述のようにmpif.h内のPARAMETER文で値が定義されていますが、この値をユーザーが関知する必要はありません。)

一方、ユーザーが何らかの理由で、図3-2-3のように例えば2つのプロセスからなるIWORLD2というグループ名(コミュニケーション)のグループを作成することもできます(作成方法は後述します)。グループを作成する際、どのプロセスをランク**0**にし、どのプロセスをランク**1**にするかをユーザーが指定します(ランク**0**と**1**は、ランク**0**~**2**とは別です)。なお、このように新規にグループを作成する事はほとんどありません。

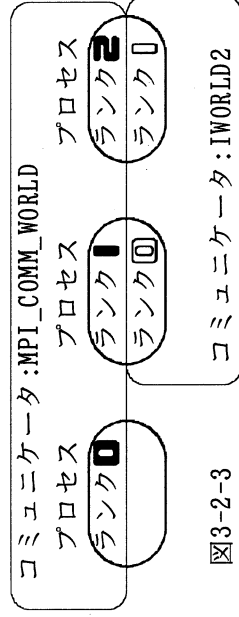


図3-2-3

上記の例を元に、図3-2-1(1)の③について説明します。③はグループに関する情報を照会するサブルーチンです。まず付録のMPI\_COMM\_SIZEの説明を一読して下さい...

上記の例では、コミュニケーションがMPI\_COMM\_WORLDであるグループに所属するプロセスの数は3なので、図3-2-1(1)の③に示すように、変数NPROCSに3が戻ります。なお本書では、全プロセス数を表す変数(整数)をNPROCSという名前に統一します。

上記のようにコミュニケーションがIWORLD2であるグループが作成されている場合、③でMPI\_COMM\_WORLDの代わりにIWORLD2を指定すると、変数NPROCSに2が戻ります。

なお、並列プログラムを1プロセスで実行すると、誤動作する可能性がありますので(4-8-2節の「プログラムの並列化方法のミス(1)」参照)、並列プログラムの開始時に、全プロセス数(NPROCS)が1以下の場合はエラーで終了させるようにして下さい(4-8-1節の「(5)並列に動作させる最低限の修正」参照)。

図3-2-1(1)の④もグループに関する情報を照会するサブルーチンです。付録のMPI\_COMM\_RANKの説明を一読して下さい...

④に示すように、ランク**0**~**2**のプロセスの変数MYRANKにそれぞれ**0**~**2**が戻ります。なお本書では、ランクを表す変数(整数)をMYRANKという名前に統一します。

上記のようにコミュニケーションがIWORLD2であるグループが作成されている場合、④でMPI\_COMM\_WORLDの代わりにIWORLD2を指定すると、変数MYRANKに**0**または**1**が戻ります。

③と④で得られたNPROCSとMYRANKの値を⑤で書き出したため、図3-2-1(2)の結果が得られました。

## ■ まとめ

以上をまとめます。まず、並列プログラムでは図3-2-1(1)の②と⑥は必須です。③と④は全プロセス数(NPROCS)あるいは自分のランク(MYRANK)を取得したい場合に使用しますが、これらもほとんどの並列プログラムでは必須です。従って、②~④はいわば三種の神器、あるいはオマジナイとして、プログラムの最初に1箇所にまとめ一度だけ実行するのがよいでしょう。そしてMPIのサブルーチンをコールしたユーザールーチンでは、①のインクルード文を指定するのを忘れないようにして下さい。

本節では、環境管理に関するサブルーチンを4つ紹介しました。それ以外に「MPI\_ABORT」というサブルーチンも知っておいた方が便利です。MPI\_ABORTは4-4-1節と4-8-1節で説明します。



■ 同じプログラムで各プロセスに異なる動作をさせるには

2章で説明した図2-1-1とその実行結果を以下に再掲します。このプログラムを並列に実行すると全プロセスが同じメッセージを書き出しました。SPMDモデルではプログラムは1つしかありませんが、図2-1-1では各プロセスの動作自体も全く同一であったため、全プロセスが同じ結果を出力したわけです。

なお、図2-1-1で説明したように、このプログラムはマシン環境によっては並列に動作しないことがあります。

```
PROGRAM MAIN
PRINT *, 'Hello MPI'
END
```

図2-1-1と同じ

```
[aoyama@node01]~/u/aoyama: mpirun -np 2 a.out ↵
Hello MPI
Hello MPI
```

それでは、1つのプログラムでありながら各プロセスで異なる動作をさせるにはどのような方法でしょうか？ 実は、先ほど紹介した『CALL MPI\_COMM\_RANK』によって得られる自分のプロセスのランク(MYRANK)が、1つのプログラムでありながら各プロセスで異なる動作をさせるための全ての根源となります。それではMYRANKの使用例を見てみましょう。

【例1】MYRANKをIF文の条件式に使用

図3-2-4(1)では、MYRANKをIF文の条件式の中で使用しています。各プロセスでは、図中の下線を引いた部分のIF文のみが真となり、二重線のステートメントが実行されます。このようにすれば、同じプログラムでありながらプロセスによって異なるステートメントを実行させることができます。

ランク0	ランク1	ランク2
<pre>PROGRAM MAIN INCLUDE 'mpif.h' CALL MPI_INIT(IERR) CALL MPI_COMM_SIZE &amp; (MPI_COMM_WORLD, NPROCS, IERR) CALL MPI_COMM_RANK &amp; (MPI_COMM_WORLD, MYRANK, IERR)</pre> <p>↓</p> <pre>IF (MYRANK == 0) THEN PRINT *, 'I am Rank 0' ELSEIF (MYRANK == 1) THEN PRINT *, 'I am Rank 1' ELSEIF (MYRANK == 2) THEN PRINT *, 'I am Rank 2' ENDIF CALL MPI_FINALIZE(IERR) END</pre>	<pre>PROGRAM MAIN INCLUDE 'mpif.h' CALL MPI_INIT(IERR) CALL MPI_COMM_SIZE &amp; (MPI_COMM_WORLD, NPROCS, IERR) CALL MPI_COMM_RANK &amp; (MPI_COMM_WORLD, MYRANK, IERR)</pre> <p>↓</p> <pre>IF (MYRANK == 0) THEN PRINT *, 'I am Rank 0' ELSEIF (MYRANK == 1) THEN PRINT *, 'I am Rank 1' ELSEIF (MYRANK == 2) THEN PRINT *, 'I am Rank 2' ENDIF CALL MPI_FINALIZE(IERR) END</pre>	<pre>PROGRAM MAIN INCLUDE 'mpif.h' CALL MPI_INIT(IERR) CALL MPI_COMM_SIZE &amp; (MPI_COMM_WORLD, NPROCS, IERR) CALL MPI_COMM_RANK &amp; (MPI_COMM_WORLD, MYRANK, IERR)</pre> <p>↓</p> <pre>IF (MYRANK == 0) THEN PRINT *, 'I am Rank 0' ELSEIF (MYRANK == 1) THEN PRINT *, 'I am Rank 1' ELSEIF (MYRANK == 2) THEN PRINT *, 'I am Rank 2' ENDIF CALL MPI_FINALIZE(IERR) END</pre>

図3-2-4(1)

```
[aoyama@node01]~/u/aoyama: mpirun -np 3 a.out ↵
I am Rank 0
I am Rank 1
I am Rank 2
```

図3-2-4(2)

## 【例2】MYRANKをD0ループの分割に使用した例

次に、MYRANKをD0ループの分割に使用する例について説明します。図3-2-5(1)は、配列Nに入っている9個のデータの合計を変数ISUMに求め、それを書き出すプログラムです(図3-2-5(2)参照)。各D0ループの反復は『D0 I=1,9』なので、これを3つのランク0,1,2のプロセスにそれぞれ『D0 I=1,3』、『D0 I=4,6』、『D0 I=7,9』のように割り振ることにします。

```
PROGRAM MAIN
INTEGER N(9)
D0 I = 1, 9
  N(I) = I
ENDDO
ISUM = 0
D0 I = 1, 9
  ISUM = ISUM + N(I)
ENDDO
PRINT *, 'ISUM = ', ISUM
END
```

図3-2-5(1)

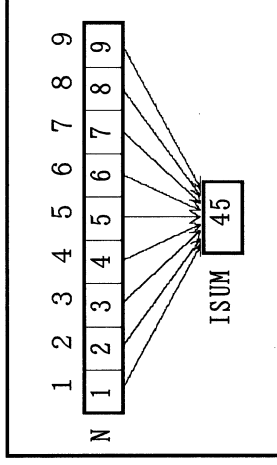


図3-2-5(2)

MYRANKを使用してD0ループを分割したプログラムを図3-2-5(3)に示します。

- まず①と②で、MYRANKの値を使用して、各プロセスが担当するISTAとIENDの値を求めます。
- 次に各D0ループの反復を『I=1,9』の代わりに『I=ISTA,IEND』とします。
- これを3プロセスで実行した場合、各プロセスは自分の担当するループ反復のみを計算し、変数ISUMには各プロセスが求めた小計が入ります。実行結果は図3-2-5(4)のようになります。

この例は実際の並列プログラムにだいたい近づいてきました。実際のプログラムでも多くの場合、この例のようにD0ループを分割して並列化を行います。ただし、この例ではD0ループの反復数(9)とプロセス数(3)があらかじめわかっており、しかもちょうど割りきれました。実際のプログラムでは、N1~N2までの反復をNPROCS個のプロセスに割り振るといように、処理量とプロセス数が任意の値でも対処できるように分割を行います。この方法については4章で説明したいと思います。

ここで1章の説明を思い出して下さい。図3-2-5(3)で行ったのは、『本当の意味の並列化とは、計算量を1/nに縮小すること(ノード数がn台の場合)』の部分です。この修正の結果、『本当の意味の並列化に伴う矛盾(副作用)』が発生しました。この例では『変数ISUMには合計でなく小計しか求まっていない』ことが矛盾(副作用)になります。そこで『矛盾(副作用)を解消するために、必要最小限、仕方なしにメッセージ交換』を行います。この例では、『小計ISUMを全プロセスで合計するために通信』を行います。

図3-2-5(3)の説明はいったんここで中断し、続きはメッセージ交換の説明を終えた3-3-6節で説明します。

ランク0

```
PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER N(9)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
```

```
①  $\overline{I}STA = \overline{MYRANK} * 3 + 1$ 
   ↳ 1
②  $\overline{I}END = \overline{I}STA + 2$ 
   ↳ 3
```

```
DO I =  $\overline{I}STA$ ,  $\overline{I}END$ 
  N(I) = I
ENDDO
ISUM = 0
DO I =  $\overline{I}STA$ ,  $\overline{I}END$ 
  ISUM = ISUM + N(I)
ENDDO
```

```
PRINT *, 'ISUM = ', ISUM
CALL MPI_FINALIZE(IERR)
END
```

ISTA IEND

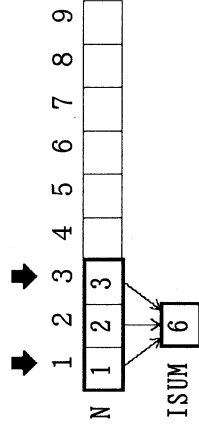


図3-2-5(3)

ランク1

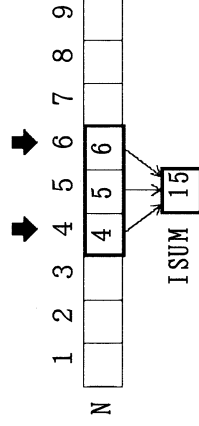
```
PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER N(9)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
```

```
①  $\overline{I}STA = \overline{MYRANK} * 3 + 1$ 
   ↳ 4
②  $\overline{I}END = \overline{I}STA + 2$ 
   ↳ 6
```

```
DO I =  $\overline{I}STA$ ,  $\overline{I}END$ 
  N(I) = I
ENDDO
ISUM = 0
DO I =  $\overline{I}STA$ ,  $\overline{I}END$ 
  ISUM = ISUM + N(I)
ENDDO
```

```
PRINT *, 'ISUM = ', ISUM
CALL MPI_FINALIZE(IERR)
END
```

ISTA IEND



ランク2

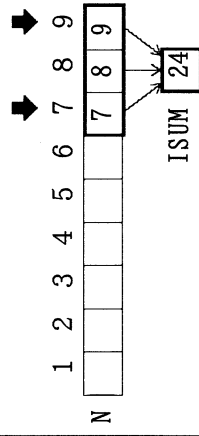
```
PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER N(9)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
```

```
①  $\overline{I}STA = \overline{MYRANK} * 3 + 1$ 
   ↳ 7
②  $\overline{I}END = \overline{I}STA + 2$ 
   ↳ 9
```

```
DO I =  $\overline{I}STA$ ,  $\overline{I}END$ 
  N(I) = I
ENDDO
ISUM = 0
DO I =  $\overline{I}STA$ ,  $\overline{I}END$ 
  ISUM = ISUM + N(I)
ENDDO
```

```
PRINT *, 'ISUM = ', ISUM
CALL MPI_FINALIZE(IERR)
END
```

ISTA IEND



[aoyama@node011]/u/aoyama: `mpirun -np 3 a.out` [ ]

```
ISUM = 6
ISUM = 15
ISUM = 24
```

図3-2-5(4)

### 3-3 集団通信

本節では、MPIで提供されているサブルーチンの中で最もよく使用する、集団通信に関するサブルーチンの使用方法について説明します。

#### 3-3-1 集団通信サブルーチンの一般的な規則

集団通信サブルーチンを使用すると、1回のサブルーチンコールで複数プロセスがメッセージ交換(データのやり取り)を一声に行うことができ、また使用方法も簡単なため大変便利です。各集団通信サブルーチンに共通する規則は次の通りです。

- その通信に参加する全プロセスが集団通信サブルーチンをコールします。
- 一声に通信を行うプロセスの範囲をコミュニケーションとしてサブルーチンコールの引数で指定します。

これを図3-3-1(1)(2)の例で示します。3つのプロセスが並列に実行しているとします。また、集団通信サブルーチンの例として、あるプロセスからその他の全プロセスにデータを送信するサブルーチン『CALL MPI\_BCAST』を取り上げます(詳細は後述します)。

図3-3-1(1)のように、引数にコミュニケーションMPI\_COMM\_WORLDを指定して『CALL MPI\_BCAST』をコールすると、コミュニケーションMPI\_COMM\_WORLDのグループに含まれるプロセス、すなわち3つの全プロセスで一斉通信を行います。これに対し、図3-3-1(2)のように、引数にコミュニケーションIWORLD2を指定してコールすると、コミュニケーションIWORLD2のグループに含まれる右の2つのプロセスのみで一斉に通信が行われます。

通常SPMDモデルでプログラムを並列化した場合、図3-3-1(2)のように一部のプロセスだけで通信を行うケースはほとんどないので、以後の説明では、全プロセスで集団通信サブルーチンを行うことを想定します。

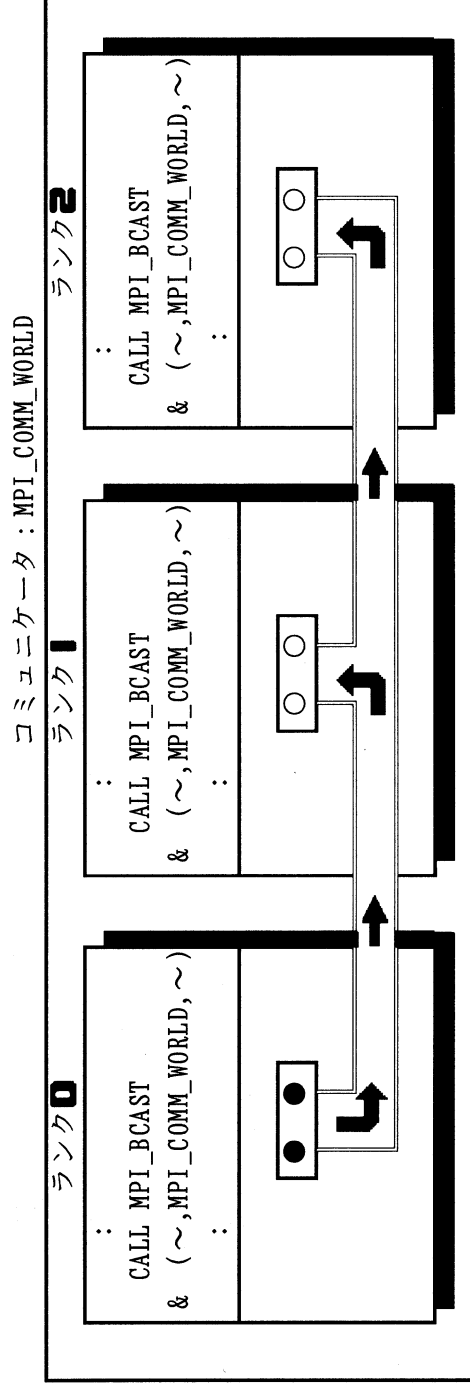


図3-3-1(1)

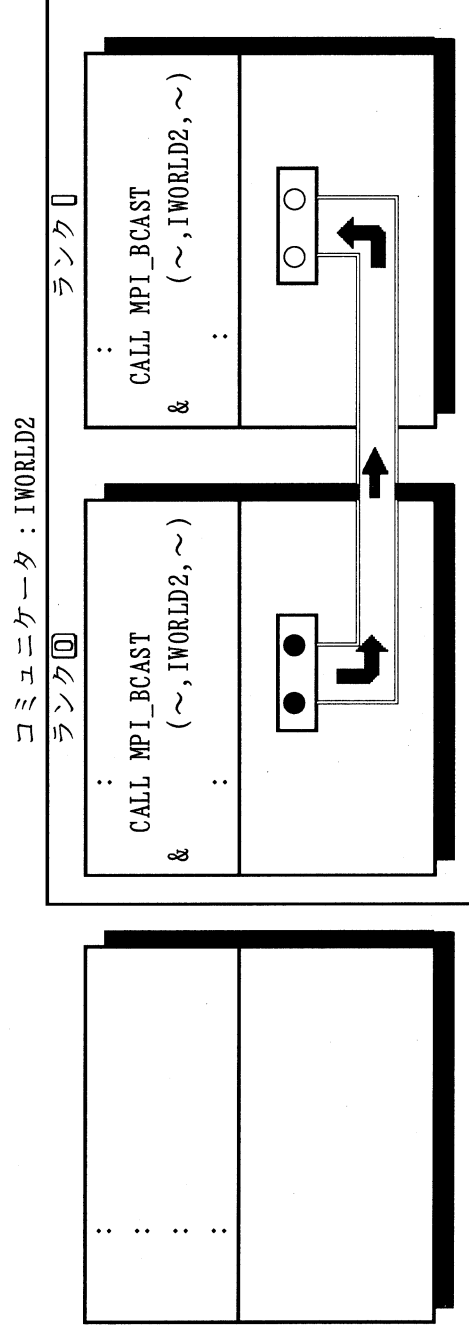


図3-3-1(2)

「並列プログラムは同期を取るのが難しい」という言葉を聞く事があります。結論を先に述べるとこれは間違いで、「並列プログラムでは、(通常)同期を取るといふことをいちいち気にする必要はない」が正解です。これについて、集団通信サブルーチンを例に以下で説明しますが、今の段階ではやや難しいので、内容は理解できなくても構いません。

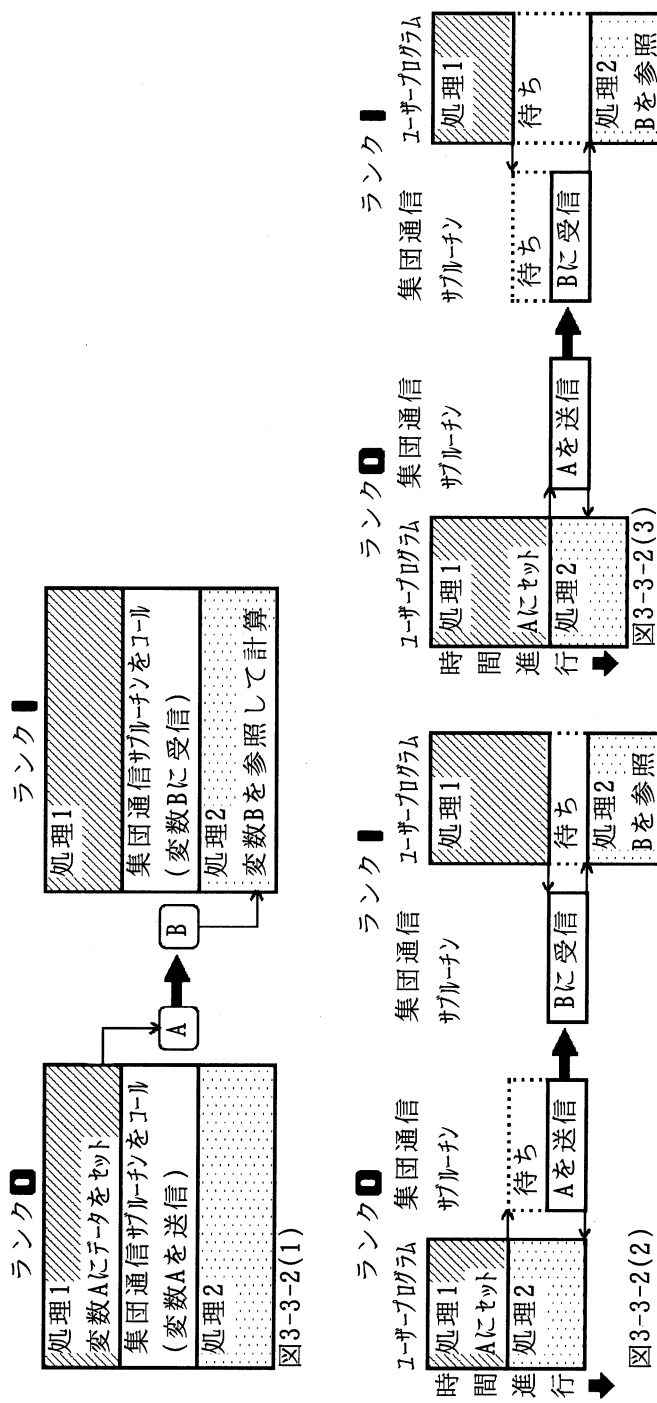
図3-3-2(1)の並列プログラムは実際には1本のプログラムですが、ランク0とランク1のプロセスで一部処理が異なる部分があるため、別々に描かれています。ランク0はまず処理1を行い、処理1の最後に変数Aにデータをセットし、集団通信サブルーチンをコールします。この集団通信サブルーチンは内部的に、変数Aをランク1に送信する処理を行うとします。ランク0は最後に処理2を実行します。

一方ランク1は処理1が終わった後、集団通信サブルーチンで変数Bを受信する処理を行います。この集団通信サブルーチンは内部的に、ランク0から送られたデータを変数Bに受信する処理を行うとします。ランク1は最後に処理2を行い、受信した変数Bを参照して計算を行います。

ランク0とランク1のプロセスは独立に動くので、どちらのプロセスが先に処理1を終了して集団通信サブルーチンをコールするかは分かりません。まず図3-3-2(2)に示すように、ランク0が先に集団通信サブルーチンをコールした場合を説明します(図の下方向は時間軸)。ランク0が集団通信サブルーチンをコールした時点で、ランク1は集団通信サブルーチンに入ります。このときランク1のユーザプログラムは、処理2に進まずに自動的に「待ち」になります(もし「待ち」にならずに先に進むと、次の処理2で、まだ受信していない変数Bを参照してしまい、誤動作するので)。そしてランク0が集団通信サブルーチンをコールした時点で通信が行われます。受信したデータが変数Bに入ると、ランク1のユーザプログラムの「待ち」は解除され、処理2に進みます。受信したデータが変数Bに入ると、ランク1のユーザプログラムの「待ち」は共に正しく動作をします。これは集団通信サブルーチンが、ユーザプログラムを先に進ませるか「待ち」にするかを適切に調整しているためです。従って集団通信サブルーチンをコールする場合、ユーザプログラム自身で、他のプロセスの処理がどこまで進んでいるのかをいちいち気にする必要はありません。言い換えると明示的に同期を取る必要はありません。

逆に図3-3-2(3)のように、ランク1が先に集団通信サブルーチンをコールした場合、ランク0はまだ処理1を行っているので、コールされた集団通信サブルーチンは「待ち」になります。このときランク1のユーザプログラムは、処理2に進まずに自動的に「待ち」になります(もし「待ち」にならずに先に進むと、次の処理2で、まだ受信していない変数Bを参照してしまい、誤動作するので)。そしてランク0が集団通信サブルーチンをコールした時点で通信が行われます。受信したデータが変数Bに入ると、ランク1のユーザプログラムの「待ち」は解除され、処理2に進みます。受信したデータが変数Bに入ると、ランク1のユーザプログラムの「待ち」は共に正しく動作をします。これは集団通信サブルーチンが、ユーザプログラムを先に進ませるか「待ち」にするかを適切に調整しているためです。従って集団通信サブルーチンをコールする場合、ユーザプログラム自身で、他のプロセスの処理がどこまで進んでいるのかをいちいち気にする必要はありません。言い換えると明示的に同期を取る必要はありません。

MPIでは、同期を取るためのメッセージ交換サブルーチンとして、MPI\_BARRIER(付録参照)が提供されています。初心者がプログラムを並列化し、なかなか答えが合わず、適当にMPI\_BARRIERを挿入したとしても答えが合うようになった、という例がたまにあります。しかし上記で説明したように、同期を取らないと答えが合わないということは通常の並列プログラムでは発生しないので、MPI\_BARRIERを挿入しないとなぜ答えが合わないのか、真の原因を把握するようにして下さい。



### 3-3-3 集団通信サブルーチンの種類

#### ■ MPIで提供されている集団通信サブルーチン

MPIで提供されている集団通信サブルーチンを下記に示します(\*はMPI-2で提供されているサブルーチンです)。各サブルーチンは以下のように分類されます。次節以降では、(グループ1)~(グループ3)のサブルーチンの代表例として、下線のサブルーチンを説明します。

なお、(その他)のサブルーチンはデータの通信自体は行ないません。

(グループ1) ◎MPI\_BCAST  
(グループ2) ○MPI\_GATHER , MPI\_SCATTER , ○MPI\_ALLGATHER , △MPI\_ALLTOALL  
○MPI\_GATHERV , MPI\_SCATTERV , ○MPI\_ALLGATHERV , △MPI\_ALLTOALLV , \*MPI\_ALLTOALLW  
(グループ3) ○MPI\_REDUCE , ○MPI\_ALLREDUCE , MPI\_SCAN , \*MPI\_EXSCAN , MPI\_REDUCE\_SCATTER  
(その他) △MPI\_BARRIER , △MPI\_OP\_CREATE , MPI\_OP\_FREE

参考までに、上記のうち使用頻度の高いサブルーチンに以下の印を付けておきます。それ以外のサブルーチンは(少なくとも私は)ほとんど使用しません。

◎:よく使用するサブルーチンです。

○:よく使用するサブルーチンですが、MPI-1では機能的に不便な点があります(3-3-5節、4-6-3節参照)。

△:たまたま使用するサブルーチンです。

### 3-3-4 MPI\_BCAST

まずMPI\_BCASTについて説明します。このサブルーチンは、1つのプロセスから、コミュニケーションで指定したグループ内の他のすべてのプロセスにメッセージを一齐に送信します。例えば図3-3-3(1)のように、大きな4の整数配列IMSGがあり、ランク0のプロセスの配列IMSGには1~4が、その他のプロセスのIMSGには0が入っているとしています。MPI\_BCASTを実行すると、ランク0のプロセスのIMSGの値がその他のプロセスのIMSGに一齐に送信されます。ここで、付録のMPI\_BCASTの説明を一読して下さい....。

図3-3-3(1)を行うプログラム例を図3-3-3(2)に、実行結果を図3-3-3(3)に示します。

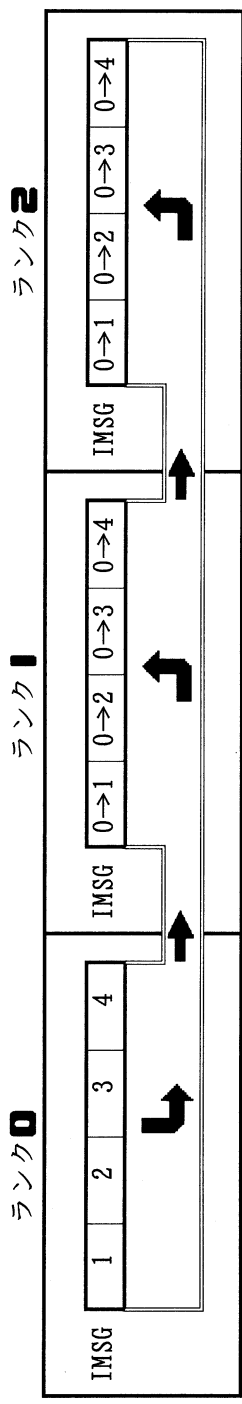


図3-3-3(1)

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER IMSG(4)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)
IF (MYRANK == 0) THEN
  DO I = 1, 4
    IMSG(I) = I
  ENDDO
ELSE
  DO I = 1, 4
    IMSG(I) = 0
  ENDDO
ENDIF
PRINT *, 'MYRANK = ', MYRANK, ' BEFORE = ', IMSG
CALL MPI_BCAST(IMG, 4, MPI_INTEGER, 0, MPI_COMM_WORLD, IERR)
PRINT *, 'MYRANK = ', MYRANK, ' AFTER = ', IMSG
CALL MPI_FINALIZE(IERR)
END
  
```

図3-3-3(2)

```

[aoyama@node01]u/aoyama: mpirun -np 3 a.out
MYRANK = 0 BEFORE = 1 2 3 4
MYRANK = 0 AFTER = 1 2 3 4
MYRANK = 1 BEFORE = 0 0 0 0
MYRANK = 1 AFTER = 1 2 3 4
MYRANK = 2 BEFORE = 0 0 0 0
MYRANK = 2 AFTER = 1 2 3 4
  
```

図3-3-3(3)

MPIが使用するインクルードファイルを指定します。  
 送信に使用するバッファを宣言します。

- ①MPIの環境を初期設定します。
- ②プロセス数NPROCSを取得します。
- ③自分のランクMYRANKを取得します。

④ランク0のプロセスだけがこの部分を実行し、配列IMSGにデータをセットします。

⑤ランク0以外のプロセスがこの部分を実行し、配列IMSGをゼロクリアします。

⑥通信前の配列IMSGの値を書き出します。

⑦ランク0からその他の全てのプロセスにIMSGの値を送信します。

⑧通信後の配列IMSGの値を書き出します。

表示される順序はタイミングによって異なります。

以下に、図3-3-3(2)の説明を行います。

●①～③を行うと、前述のように全プロセス数(NPROCS)と自分のランク(MYRANK)が得られます。ただし本プログラムではNPROCSは使用しません。次に④と⑤で、ランク0のプロセスでは配列IMSGに1～4を、その他のプロセスでは配列IMSGに0をセットします。なお、⑤では説明の都合上、配列IMSGをゼロクリアしてありますが、実際にはMPI\_BCASTによって送信されたデータによって上書きされるため、ゼロクリアする必要はありません。

●⑦で集団通信サブルーチンMPI\_BCASTをコールします。1つ目の引数(IMSG)には、通信に使用するバッファを指定します。バッファとは、送信したいデータ、または受信したデータを入れる変数や配列(の先頭アドレス)のことです。なお、本書では、バッファを送信(受信)バッファと呼ぶ場合もあります。MPI\_BCASTでは、1つ目の引数の意味が送信元のプロセスと宛先のプロセスで異なります。送信元のプロセス(本例ではランク0のプロセス)では、送信したいデータの入った送信バッファを、宛先のプロセス(本例ではランク0以外のプロセス)では、受信したデータを入れる受信バッファを指定します。⑦では全プロセスが同じ名前の配列IMSGを指定しています。

配列を指定する場合、⑦のように配列名そのもの(つまり配列の先頭アドレス)を指定する方法以外に、図3-3-3(4)のように配列の途中のアドレスを指定したり、図3-3-3(5)のようにプロセスのランクによって配列の異なるアドレスを指定してもかまいません。

⑦の前後(⑥と⑧)で全プロセスの配列IMSGの値を書き出しています。集団通信サブルーチンの使用方法を習得する段階では、このように、通信前と通信後のバッファの値をランクの値とともに書き出して、通信が行われている様子を確認するのがよいでしょう。

```

:
CALL MPI_BCAST(IMSG(2),2,MPI_INTEGER,0,MPI_COMM_WORLD,IERR)
:

```

図3-3-3(4)

```

:
CALL MPI_BCAST(IMSG(MYRANK+1),2,MPI_INTEGER,0,MPI_COMM_WORLD,IERR)
:

```

図3-3-3(5)

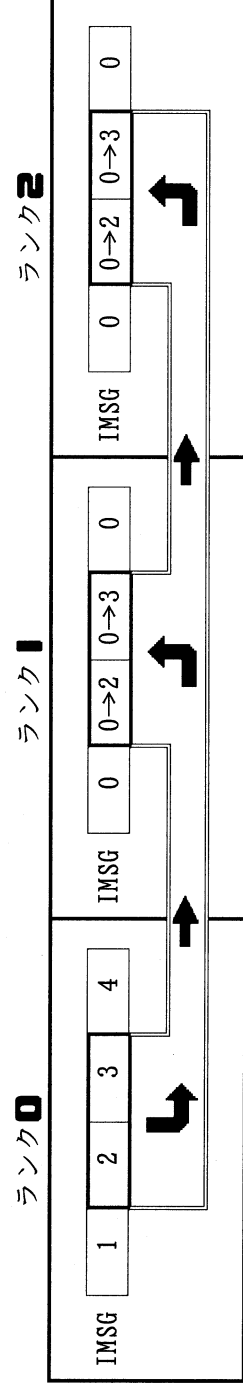


図3-3-3(6)

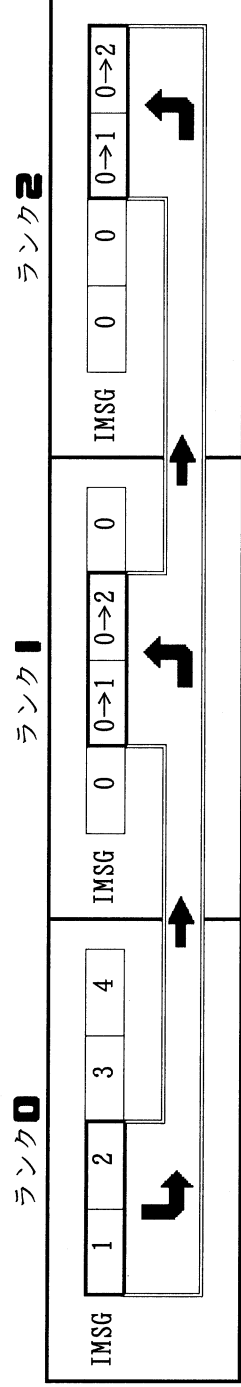


図3-3-3(7)

```

N = 4
CALL MPI_BCAST(IMSG,N,MPI_INTEGER,0,MPI_COMM_WORLD,IERR)
:

```

図3-3-3(8)



● MPIでは、送受信するデータのことをメッセージと呼びます。MPI\_BCASTの2つ目の引数には、送信するメッセージを構成するデータの個数を、3つ目の引数にはメッセージのデータ型を指定します。  
 データ型には、MPIがあらかじめ提供しているMPIの基本データ型、またはユーザーが自分で作成する派生データ型(3-5節参照)を指定します。

MPIの基本データ型とFortranのデータ型の対応を図3-3-3(9)に示します。マシン環境によっては、カッコをつけたデータ型やこれ以外のデータ型が提供されている場合があります。なお、C言語の基本データ型については3-7節を参照して下さい。

	Fortranのデータ型	MPIの基本データ型
整数型	INTEGER	MPI_INTEGER
	(INTEGER*1)	(MPI_INTEGER1)
	(INTEGER*2)	(MPI_INTEGER2)
	(INTEGER*4)	(MPI_INTEGER4)
実数型	REAL	MPI_REAL
	(REAL*2)	(MPI_REAL2)
	(REAL*4)	(MPI_REAL4)
	(REAL*8)	(MPI_REAL8)
複素数型	DOUBLE PRECISION	MPI_DOUBLE_PRECISION
	COMPLEX	MPI_COMPLEX
	(COMPLEX*16)	(MPI_COMPLEX16)
	(DOUBLE COMPLEX)	(MPI_DOUBLE_COMPLEX)
その他	LOGICAL	MPI_LOGICAL
	CHARACTER(1)	MPI_CHARACTER
	対応なし	MPI_BYTE, MPI_PACKED

図3-3-3(9)

(注)

マシン環境によっては  
 カッコをつけたデータ型などが  
 提供されている場合があります。

例えば図3-3-3(1)の例では整数型のデータを4つ送信するので、図3-3-3(2)の⑦に示すように2つ目の引数は『4』、3つ目の引数は『MPI\_INTEGER』になります。また、図3-3-3(4)(5)の例では整数型のデータを2つ送信するので、送信結果はそれぞれ図3-3-3(6)(7)のようになります。データの個数は、図3-3-3(8)の下線部のように変数で指定しても勿論かまいません。

いくつかが補足します。

- 1回の通信で送信できるバイト数はマシン環境によって異なり、無制限の場合と、制限のある場合があります。制限のある場合に、制限を越えたバイト数を送ってもエラーにならず、不完全なメッセージが送られてしまうマシン環境もありますので、お使いのマシンのマニュアルを確認して下さい。

- 集団通信ルーチンでは、送信バッファで指定するバイト数(データ個数×データ型のバイト数)と、受信バッファで指定するバイト数は一致していなければなりません。

- メッセージを、MPIのようにデータの個数とデータ型の2つで指定するのではなく、単にバイト数で指定した方が引数が1つなので簡単ですが、MPIでは個数とデータ型を分離したことによって、ユーザーが自分独自のデータ型(3-5節の派生データ型)を作成することができます。

● MPI\_BCASTの4つ目の引数には、送信元プロセスのランクを指定します。図3-3-3(1)ではランク0のプロセスが送信元なので『0』を指定します。この引数は、全プロセスが同じ値を指定する必要があります。SPMDモデルの場合、通常は図3-3-3(2)のように全プロセスが同一プログラムの同一の『CALL MPI\_BCAST』を実行します。この場合、MPI\_BCASTの4つ目の引数が④のように定数になっていれば、全プロセスが必然的に同じ値を指定することになるため問題はありません。しかし、図3-3-3(10)のような使い方は誤りです。この例では4つ目の引数にMYRANKを使用していますが、MYRANKの値はプロセスのランクによって異なっているため、どのプロセスが送信元なのかかわからなくなってしまいます。

なお、全プロセスが同じ値を指定すべき引数は、付録の各サブルーチンの引数の説明の所に明記しました。

```
CALL MPI_BCAST(IMG, 4, MPI_INTEGER, MYRANK, MPI_COMM_WORLD, IERR)
```

図3-3-3(10) ✖ 送信元プロセスのランクの誤った使用例

● MPI\_BCASTの5つ目の引数には、この通信に参加する全プロセスから構成されるグループのコミュニケータを指定します。図3-3-3(2)では前述のMPI\_COMM\_WORLDを指定しています。前にも述べましたが、5つ目の引数に指定したコミュニケータのグループに含まれる全プロセスが『CALL MPI\_BCAST』を実行する必要があります。SPMDモデルの場合、図3-3-3(2)のように、全プロセスが同一プログラムの同一の『CALL MPI\_BCAST』を実行するので通常問題はありません。

● 今までは全て、各プロセスが全く同一の『CALL MPI\_BCAST』を実行する例を示しました。SPMDモデルではこのようにプログラミングするのが一般的ですが、図3-3-3(11)のように、プロセスによって異なるステートメントの『CALL MPI\_BCAST』を実行してもかまいません(ただしこのような使用法はまれですが)。

なお、図3-3-3(11)では、下線部に示すように、プロセスによって異なる配列(ランク0のプロセスではIMG、それ以外のプロセスではKMSG)を送受信バッファとして使用していますが、このような使い方も問題ありません。図3-3-3(11)の実行結果を図3-3-3(12)に示します。

```

:
IF (MYRANK == 0) THEN
CALL MPI_BCAST(IMG, 4, MPI_INTEGER, 0, MPI_COMM_WORLD, IERR)
ELSE
CALL MPI_BCAST(KMSG, 4, MPI_INTEGER, 0, MPI_COMM_WORLD, IERR)
ENDIF
:

```

図3-3-3(11)

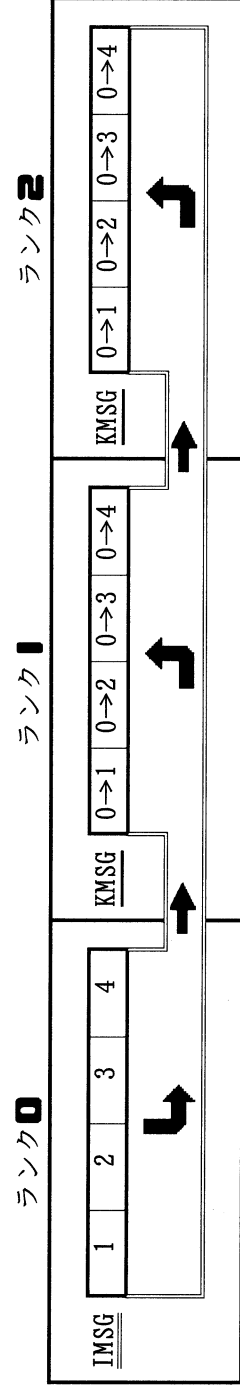


図3-3-3(12)

● MPIのサブルーチンの引数が少なかつたり指定がおかしくても、一般にコンパイル・リンクではエラーになりません。この場合、実行時にSegmentation faultなどのエラーで異常終了しますが、エラーメッセージを見ても原因が分からない場合がありますので、異常が起きた場合にはMPIのサブルーチンの引数の数や属性に誤りがないかどうかをチェックして下さい。

● MPIの仕様書(参考文献[15])では、ユーザーがMPIの使い方を間違えた場合、実行時にMPIが行うエラー処理について一般に規定していません。このため、例えばMPI\_BCASTなどの集団通信ルーチンで、送信バッファの大きさより受信バッファの大きさを誤って小さくしてしまったような場合でも、受信時にエラーにならず、エラーのリターンコード(ierr)にも正常終了の「MPI\_SUCCESS」が返るマシン環境もあります。

次に、MPI\_GATHERについて説明します。MPI\_GATHERは図3-3-4(1)のように、各プロセスの送信バッファのデータを、例えばランク0のプロセスの受信バッファに集めます。実際のプログラムでは、並列化によって各プロセスが計算した部分的な計算結果を最後にランク0のプロセスに集め、ランク0のプロセスが代表して書き出しを行うような場合に使用します。

ここで、付録のMPI\_GATHERの説明を一読して下さい。。。

以下にMPI\_GATHERに関する注意点を述べます。

- 前述のMPI\_BCASTで説明した項目はほとんどMPI\_GATHERでも適用されますが、1つだけ異なる点があります。MPI\_BCASTでは引数で指定する送受信バッファは1つだけで、送信元のプロセスでは送信バッファとして、宛先のプロセスでは受信バッファとして使用しました。実は、引数でバッファを1つだけ指定するのはMPI\_BCASTのみであり、その他の全ての集団通信サブルーチン(MPI\_BARRIERを除く)では、送信バッファと受信バッファの2つを指定します。

- 以下に示すように、(グループ2)のサブルーチンは、末尾に『V』が付いていないルーチンとついているルーチンがあります。以下で両者の相違を説明します。なお、この相違はMPI\_SCATTER(V), MPI\_ALLGATHER(V), MPI\_ALLTOALL(V), MPI\_ALLTOALL(W)でも同様です。

(グループ2) MPI\_GATHER, MPI\_SCATTER, MPI\_ALLGATHER, MPI\_ALLTOALL  
MPI\_GATHERV, MPI\_SCATTERV, MPI\_ALLGATHERV, MPI\_ALLTOALLY, MPI\_ALLTOALLV

【末尾にVのついていないルーチン(MPI\_GATHERなど)】

- (1) 図3-3-4(1)の➡に示すように、送信するメッセージの長さは全プロセスで同一でなければなりません。
- (2) 受信したメッセージは、送信元のランクの小さい順に受信バッファに入ります。

【末尾にVのついていないルーチン(MPI\_GATHERVなど)】

- (1) 図3-3-4(2)に示すように、各プロセスで異なる長さの送信メッセージを送る事ができます。
- (2) 受信したメッセージは、受信バッファの任意の位置に置くことができます。

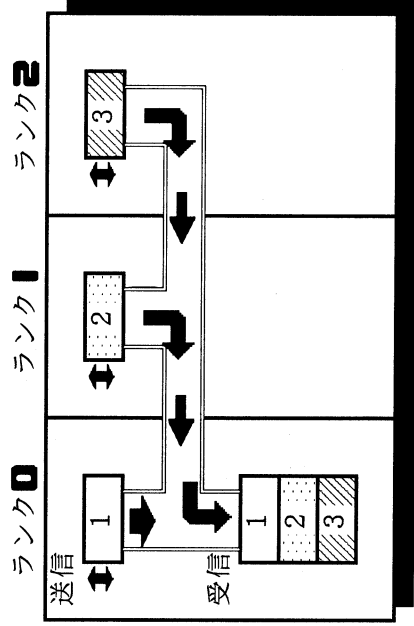


図3-3-4(1) MPI\_GATHER

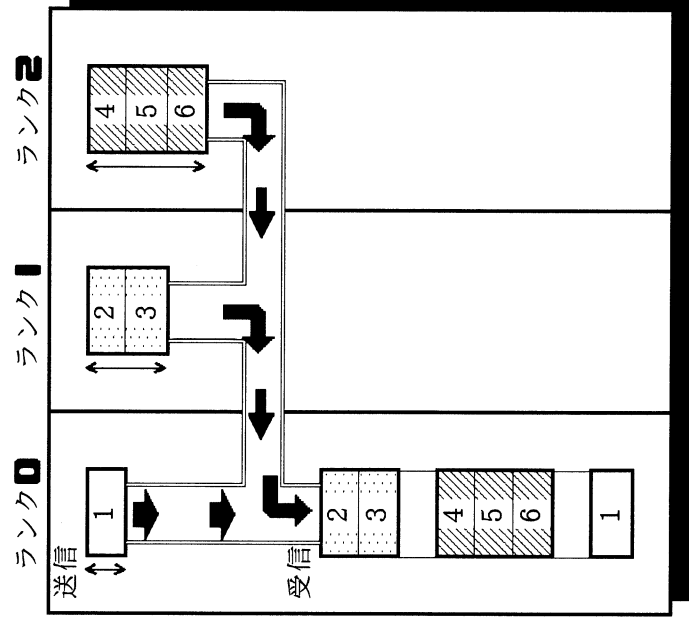


図3-3-4(2) MPI\_GATHERV

●以下に示す、送信バッファと受信バッファの両方を指定する集団通信サブルーチンでは、以下の制限があります。この制限に違反すると、タイミングによって実行結果が変わったりする可能性がありますので十分注意して下さい。

- (グループ2) MPI\_GATHER , MPI\_SCATTER , MPI\_ALLGATHER , MPI\_ALLTOALL  
MPI\_GATHERV , MPI\_SCATTERV , MPI\_ALLGATHERV , MPI\_ALLTOALLV (W)
- (グループ3) MPI\_REDUCE , MPI\_ALLREDUCE , MPI\_SCAN , MPI\_REDUCE\_SCATTER

送信バッファと受信バッファ（の実際に使用する部分）はメモリー上で重なってはならない。

実際にプログラムを並列化する場合、手抜きで、配列を縮小(4-5-7節参照)せずに並列化することがよくありますが、このとき上記の制限が問題になります(1-2節の図1-2-2(2)がその例です)。

例えば図3-3-5(1)のように、元の単体プログラムで使用していた配列Aをそのまま並列プログラムで全プロセスが使用し、各プロセスは配列Aのうち自分が計算した部分のみに値を入れたとします(図3-3-5(1)ではランク0,1,2はそれぞれA(1),A(2),A(3))。そして、これらの計算結果を最後にランク0のプロセスに集めるとします。このとき、ランク0のプロセスの配列Aは、自分が計算したA(1)以外の部分が空いているため、ランク1からの結果をA(2)に、ランク2からの結果をA(3)に入れるのが自然です。

ところがMPI\_GATHERを使用した場合、図3-3-5(1)に示すようにランク0のプロセスの送信バッファと受信バッファが重複しており、上記の制限に引っ掛かってしまいます。従って、MPI\_GATHERを使用する場合には、図3-3-5(2)に示すように新たに受信バッファ用の配列Bを宣言しなければなりません。

このような場合、後述するMPI\_SENDやMPI\_RECVなどの1対1通信サブルーチンを使用すると、図3-3-5(1)のようにデータを収集することができますが、これについては4-6-3節で説明します。

なお、MPI-2で追加されたMPI\_IN\_PLACEという変数を使用すると、この制限を回避できる場合もあります。詳細は3-8-1節を参照して下さい。

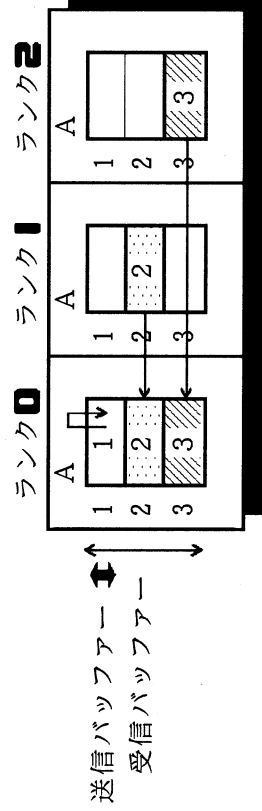


図3-3-5(1)

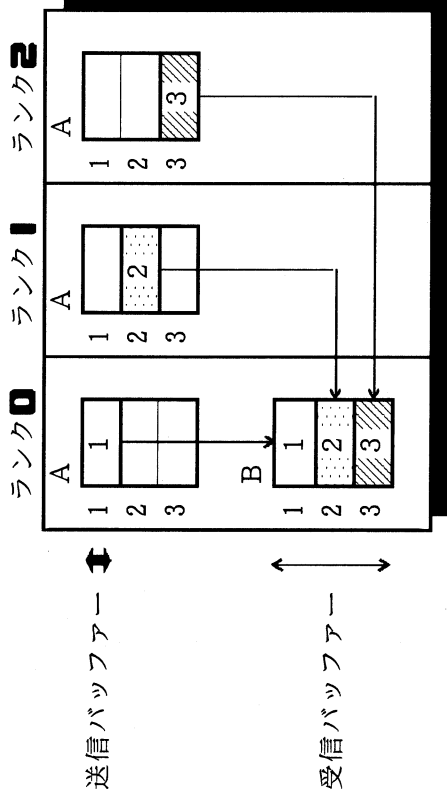


図3-3-5(2)

### 3-3-6 MPI\_REDUCE

MPI\_REDUCEは通信しながら演算を行うタイプのサブルーチンで、同じグループのサブルーチンとしては他にMPI\_ALLREDUCE, MPI\_SCAN, MPI\_REDUCE\_SCATTERがあります。まず、例によって付録のMPI\_REDUCEの説明を一読して下さい...

さて、MPI\_REDUCEでは通信しながら演算を行うわけですが、この『演算』というのは、積極的に演算しようということではなく、『通信しながら演算もいっしょに行った方が若干便利だから』と言う程度の意味合いです。とは言ってもわかりにいたいと思いますので、これを図3-3-6(1)と図3-3-6(2)で説明します。なお、この例は3-2節の図3-2-5(3)の続きです。

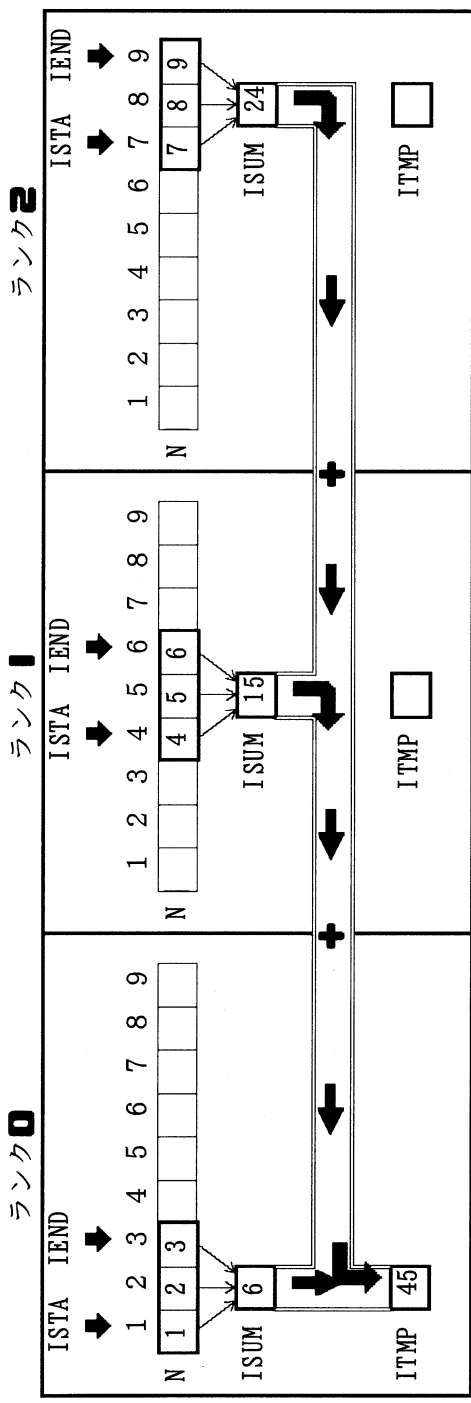


図3-3-6(1)

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER N(9)

  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  ISUM = MYRANK*3 + 1
  IEND = ISTA + 2

  DO I = ISTA, IEND
    N(I) = I
  ENDDO

  ISUM = 0
  DO I = ISTA, IEND
    ISUM = ISUM + N(I)
  ENDDO

  CALL MPI_REDUCE(ISUM, ITMP, 1, MPI_INTEGER,
    & MPI_SUM, 0, MPI_COMM_WORLD, IERR)
  ISUM = ITMP
  IF (MYRANK==0) PRINT *, 'ISUM = ', ISUM
  CALL MPI_FINALIZE(IERR)
  END
  
```

図3-3-6(2)

```

[ aoyama@node01 ] / u / aoyama : mpirun -np 3 a.out
ISUM = 45
  
```

図3-3-6(3)

MPIが使用するインクルードファイルを指定します。  
計算に使用する配列を宣言します。

CALL MPI\_INIT(IERR) MPIの環境を初期設定します。  
CALL MPI\_COMM\_SIZE(MPI\_COMM\_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。  
CALL MPI\_COMM\_RANK(MPI\_COMM\_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。

ISTA = MYRANK\*3 + 1 自分が担当する下限(ISTA)を設定します。  
IEND = ISTA + 2 自分が担当する上限(IEND)を設定します。

DO I = ISTA, IEND  
N(I) = I  
ENDDO 計算に使用する値を設定します。

ISUM = 0  
DO I = ISTA, IEND  
ISUM = ISUM + N(I)  
ENDDO ①自分が担当する部分の小計を求め、結果をISUMに入れます。

& CALL MPI\_REDUCE(ISUM, ITMP, 1, MPI\_INTEGER, ②ISUMの値を通信しながら合算し、総合計を  
MPI\_SUM, 0, MPI\_COMM\_WORLD, IERR) ランク0のプロセスのITMPに入れます。

ISUM = ITMP ③ITMPの値をISUMに戻します。  
IF (MYRANK==0) PRINT \*, 'ISUM = ', ISUM ④ISUMの値を書き出します。

CALL MPI\_FINALIZE(IERR)  
END

このプログラムでは、①が終了した直後は、図3-3-6(1)に示すように配列Nのうち各プロセスが担当した部分の小計がISUMに求まっています。最終的には小計の総合計をランク0のプロセスが持つてばよいわけですから、例えばランク1と2のプロセスのISUM(15と24)を前述のMPI\_GATHERなどでいったんランク0のプロセスに集め、その後でランク0のプロセスが『6+15+24』を計算するという方法でもよいわけです。しかしそれが若干面倒なので、通信しながら合計も同時に求めてしまおうということで、MPI\_REDUCEが提供されています。

以下、図3-3-6(2)について補足します。

●MPI\_REDUCEでは引数に通信しながら行う演算を指定します。MPIの定義済み演算(MPIがあらかじめ提供している演算)を図3-3-7(1)(2)に示します(C言語の場合は3-7節を参照して下さい)。図3-3-6(2)の②では、変数ISUMに対して加算の演算を行うため、演算には『MPI\_SUM』を指定しました。また図3-3-7(1)(2)からわかるように、各演算には、その演算で使用することのできるMPIのデータ型が決められています。なお、図3-3-7(1)(2)はマシンの環境によって若干異なることがあります。

MPIの定義済み演算の中でよく使用するのはMPI\_SUM、MPI\_MAX、MPI\_MINです。MPI\_PRODは例えば階乗の計算を並列化する場合などに使用します。MPI\_MAXLOC、MPI\_MINLOCについては3-3-6-1節で説明します。図3-3-7以外のデータ型で演算を行いたい場合(例えば倍精度複素数の加算)、3-3-6-2節で説明するように、ユーザーが自分で演算を定義することができます。

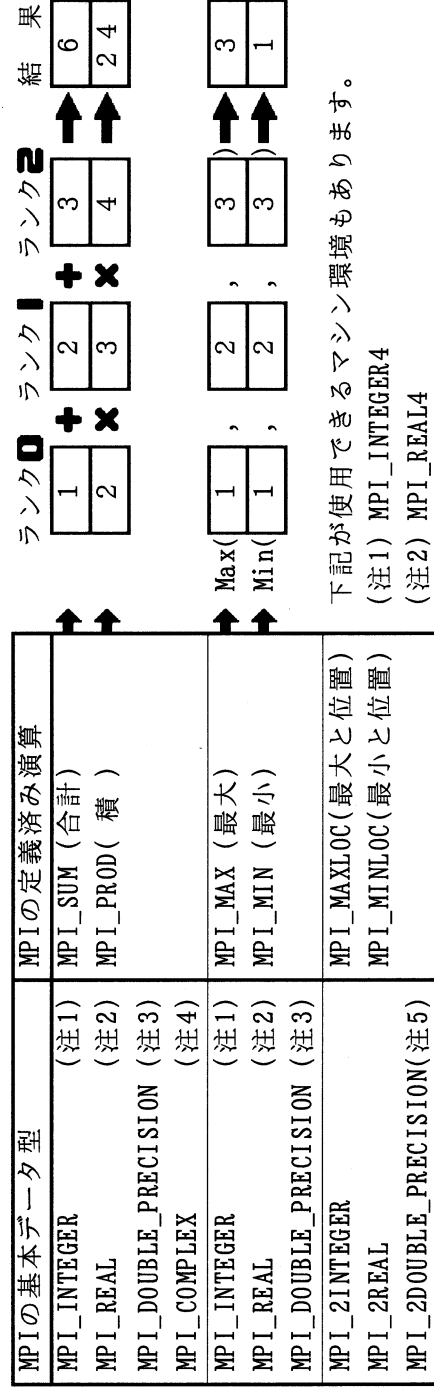


図3-3-7(1)

ただし、MPI\_2INTEGER = {MPI\_INTEGER, MPI\_INTEGER}  
MPI\_2REAL = {MPI\_REAL, MPI\_REAL}  
MPI\_2DOUBLE\_PRECISION = {MPI\_DOUBLE\_PRECISION, MPI\_DOUBLE\_PRECISION}

下記が使用できるマシンの環境もあります。

- (注1) MPI\_INTEGER4
- (注2) MPI\_REAL4
- (注3) MPI\_REAL8
- (注4) MPI\_COMPLEX16, MPI\_DOUBLE\_COMPLEX
- (注5) MPI\_2REAL8

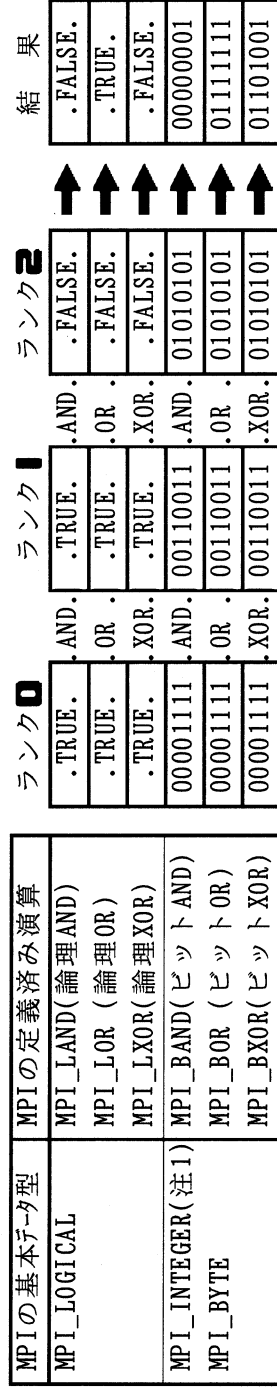


図3-3-7(2) (注1) 図3-3-7(1)の(注1)と同じ。

●D0ループの並列化と、MPI\_REDUCE(またはMPI\_ALLREDUCE)、MPIの定義済み演算の関係を整理します。

図3-3-6(2)の2つ目のループのように、並列化したループの中に合計や内積を求める計算が含まれている場合、ループを終了した直後に、「各プロセスが求めた合計の総合計」を求めるために、MPI\_SUMを使用し、MPI\_REDUCEまたはMPI\_ALLREDUCEで通信を行う必要があります。

また、並列化したループの中に、最大(最小)値を求める計算が含まれている場合、ループを終了した直後に、「各プロセスが求めた最大(最小)値の全体の最大(最小)値」を求めるために、MPI\_MAX(MPI\_MIN)を使用して、MPI\_REDUCEまたはMPI\_ALLREDUCEで通信を行う必要があります。

● MPI\_GATHERでも説明しましたが、送信バッファと受信バッファ（の実際に使用する部分）はメモリー上で重なってはいけませんので注意して下さい。なお、MPI-2で追加されたMPI\_IN\_PLACEという変数を使用すると、この制限を回避できる場合もあります。詳細は3-8-1節を参照して下さい。

図3-3-6(2)では、通信の最終結果は変数ISUMではなく受信バッファA-ITMPに入ります。これにともない④のISUMもITMPに修正する必要があります。しかし、元のプログラムの④以降で変数ISUMを何度も使用して、ISUMをITMPに直すのが面倒な場合、あるいは元のプログラムをあまり修正したくない場合には、③のようにITMPを元のISUMに代入し直し、以後はISUMをそのまま使用するという方法もあります。

● 実数の合計や内積を求めるD0ループを並列化すると、単体版と並列版で結果が若干異なったり、あるいは同じ並列版でもプロセス数によって結果が異なる場合があります。その理由を以下に示します。

a, b, cが実数のとき、数学の世界では下記の式が成立します。

$$(a + b) + c = a + (b + c)$$

一方計算機の世界では、取り扱う桁数が有限なので、実数の加算の計算順序を変えると、丸め誤差や桁落ちの影響で計算結果が若干変わることがあり、厳密には下記のようなになります。

$$(a + b) + c \neq a + (b + c)$$

例えば配列A(1)～A(8)の合計を求めるプログラムを単体版、並列版(2プロセス)、並列版(4プロセス)で行った場合、以下の(A), (B), (C)のように加算が行われます。[ ]を先に計算するため、それぞれ計算順序は異なっており、その結果、合計が若干変わることがあります。

$$(A) \text{ 単体版} \quad : \quad A(1) + A(2) + A(3) + A(4) + A(5) + A(6) + A(7) + A(8)$$

$$(B) \text{ 並列版(2プロセス)} : \quad \underbrace{[A(1) + A(2)] + A(3) + A(4)}_{\text{ランク0}} + \underbrace{[A(5) + A(6) + A(7) + A(8)]}_{\text{ランク1}}$$

$$(C) \text{ 並列版(4プロセス)} : \quad \underbrace{[A(1) + A(2)]}_{\text{ランク0}} + \underbrace{[A(3) + A(4)]}_{\text{ランク1}} + \underbrace{[A(5) + A(6)]}_{\text{ランク2}} + \underbrace{[A(7) + A(8)]}_{\text{ランク3}}$$

また(C)の波線部に示す各プロセス間の加算を集団通信ルーチンMPI\_REDUCE(MPI\_SUMを指定)を用いて行った場合、各プロセスの加算順序はマシン環境によって異なり、以下のいずれかになります。

[1] 各プロセスの加算順序は常に同じ

この場合、同じプロセス数であれば合計は必ず同じになります。

[2] 各プロセスの加算順序は早いもの勝ち

この場合、同じプロセス数でもタイミングによって合計が若干変わることがあります。

例えば連立一次方程式の反復解法(CG法など)では、答えが収束するまで(計算結果がある条件を満足するまで)計算を何度も繰り返します。この解法で現れる合計や内積のD0ループを並列化すると、上記の理由で、単体版、並列版(2プロセス)、並列版(4プロセス)で収束回数異なる現象が発生することがあります。

また上記[2]のマシン環境の場合は、例えば同じ並列版(4プロセス)で実行しても、収束回数が異なる現象が発生することがあります。例えばMPI\_ALLREDUCEの場合であれば、各プロセスの合計をMPI\_ALLGATHERなどで各プロセスに集めた後、各プロセスが例えばランクの小さい方から順番に加算するようにすれば、回避することができます。ただし恐らくMPI\_ALLREDUCEより速度は遅くなると思われれます。

### 3-3-6-1 MPI\_MAXLOC/MPI\_MINLOC

前述の図3-3-7(1)の演算MPI\_MAXLOCとMPI\_MINLOCの使用方法を説明します。図3-3-8(2)のD0ループでは図3-3-8(1)の配列Nの最大値(IMAX=90)と配列内の位置(ILOC=4)を求めています。このループを並列化した図3-3-8(3)では、まず各プロセスは②で、配列Nのうち自分の担当範囲(ISTA~IEND)の最大値と位置を求めます。これを④のMPI\_REDUCEで通信しながら演算し、全プロセスでの最大値とその位置を求めたいのですが、前述の演算MPI\_MAXでは最大値しか求められません。このような場合、①で大きさ2の送受信バッファIRECVを宣言し、③でISENDの1つ目と2つ目の要素に最大値(IMAX)と位置(ILOC)を代入します。

④のMPI\_MAXLOCは、最大値とそれに付随する位置を同時に求める演算で、MPI\_2INTEGERは連続した2箇の整数を意味するデータ型です。2箇の整数が1つの単位となるので、④の3つ目の引数は「1」となります。④を実行すると、図3-3-8(1)に示すように、ランク0のプロセスの受信バッファIRECVの1つ目の要素には全プロセスの最大値が、2つ目の要素にはその最大値の位置が入ります。以下に補足を述べます。

- 異なるプロセスの最大(小)値が同じで位置の値が異なる場合、位置の値は小さい方が取られます。
- MPI\_MAXLOCで指定する「位置」の中身は、「位置」の値でなくとも(任意の付加情報でも)構いません。
- MPI\_MAXLOCでは位置の情報は1つしか指定できません。2つ以上指定したい場合は3-3-6-2節の【例2】を参照して下さい。

●図3-3-8(1)では最大値と位置がともに整数だったので、送受信バッファは大きさ2の整数とし、MPIで提供する基本データ型のMPI\_2INTEGERを使用しました。最大値が実数で位置が整数の場合、C言語のMPIでは、MPIの基本データ型として構造型(MPI\_FLOAT\_INTなど)が用意されています(3-7節参照)。

一方FortranのMPIでは構造型が用意されていないので、図3-3-8(4)のように、送受信バッファは実数とし、位置の整数値を実数に変換して通信し、後で再び整数に戻すのが1つの方法です。

この場合、MPIの基本データ型は、図3-3-7(1)に示すように、送受信バッファは単精度であればMPI\_2REAL、倍精度であればMPI\_2DOUBLE\_PRECISIONとなります。

ただし単精度実数の場合、整数の方が単精度実数よりも有効桁数が多いため、位置の値が大きい場合は(例えば66666666)、右の例に示すように、整数値を単精度実数に代入した時に、元の値と違う値になってしまうので注意して下さい。

別の方法として、構造型の送受信バッファ(最大値が実数、位置が整数)を使用する演算をユーザーが自分で定義する方法を、

3-3-6-2節の【例1】に示します。

(結果) 666666666 666666664.0 666666664

```
INTEGER I, J (整数)
REAL A (単精度実数)
I = 666666666
A = I
J = A
PRINT *, I, A, J
```

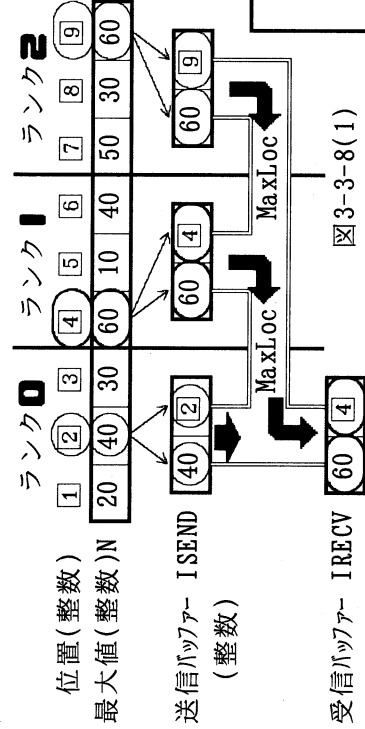


図3-3-8(1)

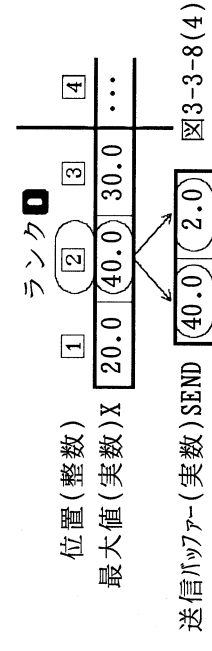


図3-3-8(4)

```
INTEGER N(9), ISEND(2), IRECV(2)
:
IMAX = -999
DO I = 1, 9
  IF (N(I) > IMAX) THEN
    IMAX = N(I)
    ILOC = I
  ENDIF
ENDDO
ISEND(1) = IMAX
ISEND(2) = ILOC
CALL MPI_REDUCE(ISEND, IRECV, 1, MPI_2INTEGER,
& MPI_MAXLOC, 0, MPI_COMM_WORLD, IERR)
PRINT *, 'MAX=' , IRECV(1), 'LOCATION=' , IRECV(2)
```

図3-3-8(3)

```
INTEGER N(9)
:
IMAX = -999
DO I = 1, 9
  IF (N(I) > IMAX) THEN
    IMAX = N(I)
    ILOC = I
  ENDIF
ENDDO
PRINT *, 'MAX=' , IMAX, 'LOCATION=' , ILOC
```

図3-3-8(2)



### 3-3-6-2 ユーザーが定義する演算

前述の図3-3-7(1)(2)と同様の「通信しながら行う演算」をユーザーが自分で作成することができます。どういうときに作成が必要になるかについては本節の後半で説明します。

以下で「通信しながら行う演算」の作成方法を説明します。以下の説明で、ユーザーが作成する通信しながら行う演算を⊕(例えば加算)で表します。また3プロセスで実行する場合を想定し、ランク**0,1,2**の各プロセスが提供する送信データを**0,1,2**で表すことにします。

#### ■ 結合則と交換則

##### 【結合則】

演算⊕は結合則を満足する必要があります。すなわち  $(0 \oplus 1) \oplus 2 = 0 \oplus (1 \oplus 2)$  である必要があります。ここで  $(0 \oplus 1) \oplus 2$  は、まずランク**0**とランク**1**の演算が行われ、次にその中間結果とランク**2**の演算が行われることを意味します。

##### 【交換則】

#### [演算⊕が交換則を満足する場合]

演算⊕が、結合則の他に交換則を満足する場合  $(0 \oplus 1) = 1 \oplus 0$  など、後述する図3-3-9の②の2つ目の引数に「.TRUE.」を指定します。これを指定すると、通信しながら行う演算の順序は以下のいずれかになります。言いかえると、以下の中に正しくない演算順序がある場合、「.TRUE.」を指定することはできません。実行中に、以下などの順序で演算が行われるかはマシン環境によって異なります。一般にそのマシン環境で高速になるような順序で行われます。

**0**⊕(**1**⊕**2**) (**1**⊕**2**)⊕**0**    **1**⊕(**0**⊕**2**) (**0**⊕**2**)⊕**1**    **2**⊕(**0**⊕**1**) (**0**⊕**1**)⊕**2**  
**0**⊕(**2**⊕**1**) (**2**⊕**1**)⊕**0**    **1**⊕(**2**⊕**0**) (**2**⊕**0**)⊕**1**    **2**⊕(**1**⊕**0**) (**1**⊕**0**)⊕**2**

#### [演算⊕が交換則を満足しない場合]

演算⊕が、結合則のみ満足し交換則を満足しない場合  $(0 \oplus 1) \neq 1 \oplus 0$  など、後述する図3-3-9の②の2つ目の引数に「.FALSE.」を指定します。これを指定すると、演算はランクの小さい方から昇順に、以下のいずれかの順序で行われます。例えば3つの行列A, B, Cの乗算の場合、 $(A \cdot B) \cdot C = A \cdot (B \cdot C)$  で結合則は満足しますが、 $A \cdot B \neq B \cdot A$  で交換則は満足しないので、こちらに該当します。

**0**⊕(**1**⊕**2**) (**0**⊕**1**)⊕**2**

#### ■ 例(通信しながら整数の加算を行う演算)

例として、「通信しながら整数の加算を行う演算」(図3-3-7(1)のMPI\_SUMに相当)をユーザーが作成する方法を、図3-3-9のプログラムで説明します。

- ②で、ユーザーが作成する演算をMPIに登録するサブルーチン『MPI\_OP\_CREATE』(詳細は付録参照)を一度だけコールします。
  - 1つ目の引数では、通信しながら行う演算のロジックを記述する④のサブルーチン名を指定します。また①で④のサブルーチン名をEXTERNAL文で外部関数として宣言します。
  - 2つ目の引数では、前述の[演算⊕が交換則を満足する場合は「.TRUE.」を、[演算⊕が交換則を満足しない場合は「.FALSE.」を指定します。
  - 3つ目の引数では、③の通信で使用する演算の名前(図3-3-7(1)のMPI\_SUMに相当する任意の名前の整数)を指定します。
- ②を実行すると、演算がMPIに登録され、そのID(識別子)が変数ISUMに戻ります。この変数を③の5つ目の引数で指定します。なお、②を含むサブルーチン以外のサブルーチンで③の通信を行う場合は、変数ISUMをそのサブルーチンに引数、COMMON文、MODULE文などで渡す必要があります。
- ③の通信を実行すると、内部的に④がコールされ、引数IN, INOUT, LEN, ITYPE(名前は任意)に値が渡されます(引数の詳細は後述します)。なおC言語の場合、④の部分のANSI-Cプロトタイプは以下のようにになります(詳細は参考文献[15][27]の4.9.4節参照)。

```
typedef void 関数名( void *invec, void *inoutvec, int *len, MPI_Datatype *datatype);
```

- 以下でサブルーチンSAMPLEの動作を説明します。

【引数が「.TRUE.」の場合】[演算が交換則を満足する場合]

②の2つ目の引数が「.TRUE.」の場合、動作は以下ようになります。

- ⑤で引数INとINOUTを、③の4つ目の引数と同じ型の変数または配列で宣言します。INとINOUTには、自分のプロセスの送信バッファ(③のISEND)で供給した値と、他のプロセスから送られた計算の中間結果が渡されますが、どちらがどちらに渡されるか(図3-3-10(1)と(2)のどちらになるか)はマシン環境によって異なります。
- 引数LEN(整数)には、③の3つ目の引数で指定した「要素数」(本例では1)が渡されます。
- 引数ITYPE(整数)には、③の4つ目の引数で指定した「データ型」(例えばMPI\_INTEGER)が渡ります。ただしFortranの場合、この値を使用することはほとんどないと思われます。
- ⑥でINとINOUTを使用して、このプロセスでの計算を要素数分だけ行い、それを図3-3-10(1)(2)に示すようにINOUTに保管します。INOUTはこのプロセスまでの中間結果として、次のプロセスに送られます。
- ④のサブルーチン内ではMPIの通信ルーチンを使用することはできません。ただしエラーが発生した時にMPI\_ABORTを使用することはできません。

【引数が「.FALSE.」の場合】[演算が交換則を満足しない場合]

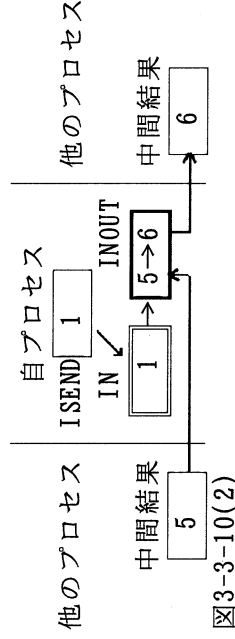
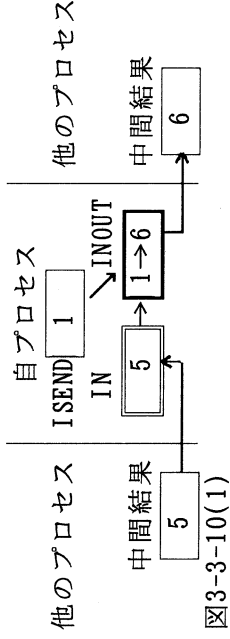
②の2つ目の引数が「.FALSE.」の場合、以下の点で「.TRUE.」の場合と異なります。自プロセスと、中間結果を送ってきた他プロセスのうち、ランクの小さい方のプロセスのデータ(ランクの小さいプロセスが他プロセスならそのプロセスから送られた中間データ、自プロセスなら送信バッファ-ISENDのデータ)がINに入り、ランクの大きいプロセスのデータがINOUTに入ります。例えば0~4の5プロセスで実行している場合、**ランク2**のプロセスでは図3-3-10(3)(4)のようになります。

演算が交換則を満足していないので、INとINOUTにどちらのプロセスのデータが入っているかを考慮して⑥の演算を行う必要があります。引数を「.FALSE.」に設定した例を、後述する【例2】で紹介します。

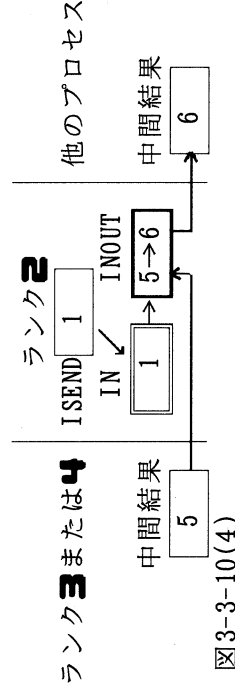
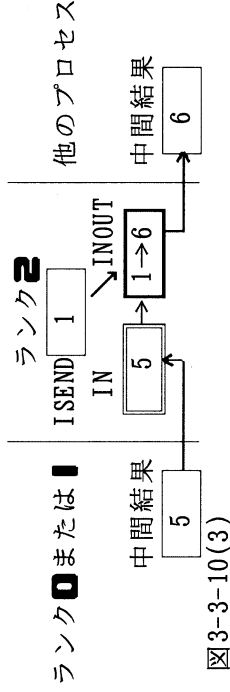
<pre>EXTERNAL SAMPLE CALL MPI_OP_CREATE(SAMPLE, .TRUE., ISUM, IERR) 送信バッファ-ISENDに送信データをセットする。 CALL MPI_REDUCE(ISEND, IRECV, 1, MPI_INTEGER, &amp; ISUM, 0, MPI_COMM_WORLD, IERR) :</pre>	<pre>SUBROUTINE SAMPLE(IN, INOUT, LEN, ITYPE) INTEGER IN(*), INOUT(*) DO I = 1, LEN   INOUT(I) = IN(I) + INOUT(I) ENDDO END</pre>
--	---

図3-3-9

【引数が「.TRUE.」の場合】(以下のいずれか)



【引数が「.FALSE.」の場合】



多くの場合、MPIの定義済み演算だけで事が足りりますが、ユーザーが「通信しながら行う演算」を作成しなければならぬ例をいくつか紹介します。

**【例1】最大(小)値が実数、位置が整数**

3-3-6-1節で説明したように、Fortranの場合、最大(小)値が実数で位置が整数という送受信データを、MPI\_REDUCEなどでMPI\_MAXLOCなどを使用して通信することができます。このデータが通信できるように、MPI\_MAXLOCなどに相当する演算をユーザーが作成する方法を図3-3-11(2)のプログラムで説明します。

- 図3-3-11(1)に示すような、単精度実数Aと整数Iから構成される構造体KOUZOUTAIを図3-3-11(2)の①で定義し、この構造体型の送信バッファA-Xと受信バッファA-XXを④で宣言します。
- ⑨のサブルーチンPARAMAXLOCで、通信しながら行うMPI\_MAXLOCに相当する演算を定義します。⑩で最大値が大きい方のプロセスの最大値と位置の値を、⑪でXINOUTに入れて次のプロセスに渡します。
- 図3-3-11(1)のように異なるプロセスの最大値(⑨と⑩)が同じで位置の値(⑨と⑩)が異なる場合、MPI\_MAXLOCと同様に、位置の値は小さい方(⑨)を取るようにします。このため、⑫で異なるプロセスの最大値が同じ場合、⑬で位置の値は小さい方を取るようにします。
- ⑤でサブルーチンPARAMAXLOCを、通信しながら行う演算IMAXLOCとしてMPIに登録します。
- ③と⑥で、この構造体を表す派生データ型(3-5-4節参照)IKOUZOUを作成します。IKOUZOUは⑮の通信で使用します。
- ⑦で、構造体の送信バッファA-Xに、通信したい最大値(単精度実数)と位置(整数)を代入します。
- ⑧で、新たに作成した演算IMAXLOCを使用してMPI\_REDUCEの通信を行います。すると図3-3-11(1)に示すように、全プロセスの実数の最大値(本例では60.0)とそれに付随する整数の位置(本例では⑨と⑩のうち小さい方の⑨)がランク0の受信バッファA-XXに入ります。

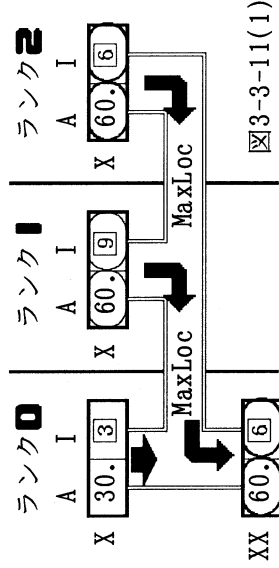


図3-3-11(1)

```

MODULE PARA
  TYPE :: KOUZOUTAI
    REAL A
    INTEGER I
  END TYPE KOUZOUTAI
END

PROGRAM MAIN
  USE PARA
  INCLUDE 'mpif.h'
  EXTERNAL PARAMAXLOC
  INTEGER IBLOCK(2), IDISP(2), ITYPE(2)
  TYPE(KOUZOUTAI) X, XX
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE
  & (MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK
  & (MPI_COMM_WORLD, MYRANK, IERR)
  CALL MPI_OP_CREATE
  & (PARAMAXLOC, .TRUE., IMAXLOC, IERR)
  
```

```

IBLOCK(1) = 1
IBLOCK(2) = 1
IDISP(1) = 0
IDISP(2) = 4
ITYPE(1) = MPI_REAL
ITYPE(2) = MPI_INTEGER
CALL MPI_TYPE_STRUCT
& (2, IBLOCK, IDISP, ITYPE, IKOUZOU, IERR)
CALL MPI_TYPE_COMMIT(IKOUZOU, IERR)
:
X%A = ~ 最大値(単精度実数)
X%I = ~ 位置(整数)
CALL MPI_REDUCE(X, XX, 1, IKOUZOU, IMAXLOC,
& 0, MPI_COMM_WORLD, IERR)
:
SUBROUTINE PARAMAXLOC(XIN, XINOUT, LEN, ITYPE)
  USE PARA
  TYPE(KOUZOUTAI) XIN(*), XINOUT(*)
  DO I=1, LEN
    IF (XINOUT(I)%A < XIN(I)%A) THEN
      XINOUT(I)%A = XIN(I)%A
      XINOUT(I)%I = XIN(I)%I
    ELSEIF (XINOUT(I)%A == XIN(I)%A) THEN
      XINOUT(I)%I = MIN(XINOUT(I)%I, XIN(I)%I)
    ENDIF
  ENDDO
END
  
```

図3-3-11(2)

【例2】最大値(最小値)と2つの位置

図3-3-12(1)に示す2次元配列Zの最大値(本例では60.0)と、その1次元目と2次元目の位置を求めるループを並列に実行した場合、ループ終了後に、各プロセスが求めた最大値の最大値とその2つの位置を求めるため、MPI\_REDUCEなどの通信を行います。ところが3-3-6-1節で説明したように、MPI\_MAXLOCでは、最大値に付随する位置は1つしか指定できません。2つ以上指定したい場合は、MPI\_MAXLOCに相当する演算をユーザが定義する必要があります。プログラム例を図3-3-12(3)に示しますが、図3-3-11(2)とほとんど同じなので、異なる部分のみ説明します。

- ①の構造体にJを追加し、構造体の派生データ型を作成する③の部分とします(3-5-4節参照)。
- ④でJ方向の位置を設定します。
- 本例のように、異なるプロセスの最大値(60.)が同じで位置の値が異なる場合、図3-3-11(2)ではMPI\_MAXLOCと同様に位置の値は小さい方を取りましたが、本例では位置が2つあってこの方法は適用できないため、ランクの小さいプロセスの値を取ります。まず②の2つ目の引数を「.FALSE.」とします。これによって、本節の前半で説明したように、⑤でランクの小さい方の位置が配列XINに、ランクの大きい方の位置が配列XINOUTに入ります。⑥で2つのプロセスの最大値が等しいときは、⑦でランクの小さいプロセスの位置を採用するので、上記下線部が実現されます。

J →	ランク0	ランク1	ランク2	ランク3	ランク4	ランク5	ランク6	ランク7	ランク8	ランク9
I ↓ ①	10.	20.	10.	30.	20.	40.	10.	20.	60.	60.
②	10.	10.	20.	10.	20.	10.	40.	20.	30.	
③	20.	30.	10.	20.	10.	30.	30.	10.	20.	
④	20.	20.	20.	30.	60.	40.	20.	40.	10.	
⑤	10.	10.	10.	40.	30.	10.	30.	20.	10.	

図3-3-12(1)

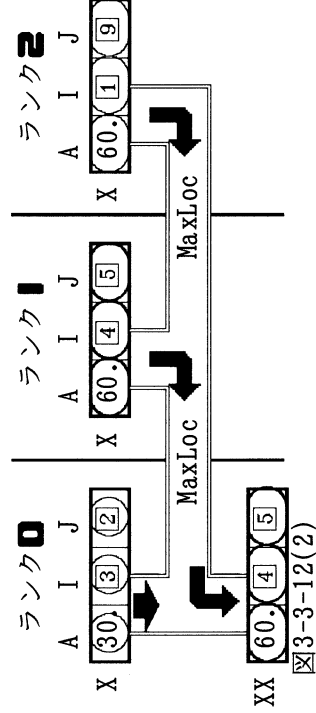


図3-3-12(2)

```

MODULE PARA
TYPE::KOUZOUTAI
  REAL A 最大値(単精度実数)
  INTEGER I I方向の位置(整数)
  INTEGER J J方向の位置(整数)
END TYPE KOUZOUTAI
END

PROGRAM MAIN
USE PARA
INCLUDE 'mpif.h'
EXTERNAL PARA_MAXLOC2
INTEGER IBLOCK(2), IDISP(2), ITYPE(2)
TYPE(KOUZOUTAI) X, XX
REAL Z(5,9)

CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
CALL MPI_OP_CREATE
& (PARA_MAXLOC2, .FALSE., IMAXLOC2, IERR) ②
IBLOCK(1) = 1
IBLOCK(2) = 2
IDISP(1) = 0
IDISP(2) = 4
  
```

```

ITYPE(1) = MPI_REAL
ITYPE(2) = MPI_INTEGER
CALL MPI_TYPE_STRUCT
& (2, IBLOCK, IDISP, ITYPE, IKOUZOU, IERR)
CALL MPI_TYPE_COMMIT(IKOUZOU, IERR)
:
X%A = ~ 最大値(単精度実数)
X%I = ~ I方向の位置(整数)
X%J = ~ J方向の位置(整数) ④
CALL MPI_REDUCE(X, XX, 1, IKOUZOU, IMAXLOC2,
& 0, MPI_COMM_WORLD, IERR)
:
SUBROUTINE PARA_MAXLOC2(XIN, XINOUT, LEN, ITYPE)
USE PARA
TYPE(KOUZOUTAI) XIN(*), XINOUT(*)
DO I=1, LEN
  IF (XINOUT(I)%A <= XIN(I)%A) THEN ⑥
    XINOUT(I)%A = XIN(I)%A
    XINOUT(I)%I = XIN(I)%I
    XINOUT(I)%J = XIN(I)%J
  ENDIF
ENDDO
END
  
```

図3-3-12(3)

【例3】倍精度複素数の加算

前述の図3-3-7(1)に示すように、MPIの仕様書(参考文献[15],[27])には、MPI\_SUM、MPI\_PRODの演算に倍精度複素数(MPI\_COMPLEX16またはMPI\_DOUBLE\_COMPLEX)が含まれていません。実際には多くのマシン環境では使うことができますが、もし使えない場合、ユーザーが「通信しながら倍精度複素数の加算を行う演算」を定義する必要があります。データの動きとプログラム例を図3-3-13(1)(2)に示します。

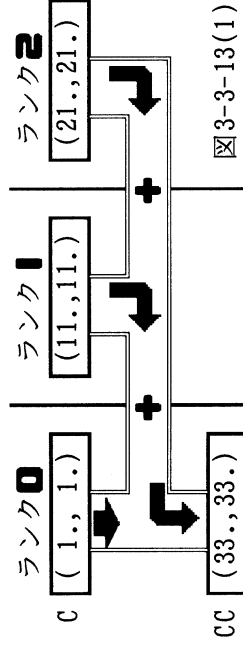


図3-3-13(1)

```

CALL MPI_OP_CREATE
& (PARA_SUMC16,.TRUE., ISUMC16, IERR)
C = 複素数の値を設定
CALL MPI_REDUCE(C, CC, 1, MPI_DOUBLE_COMPLEX,
& ISUMC16, 0, MPI_COMM_WORLD, IERR)
:
SUBROUTINE PARA_SUMC16(CIN, CINOUT, LEN, ITYPE)
COMPLEX*16 CIN(*), CINOUT(*)
DO I = 1, LEN
  CINOUT(I) = CIN(I) + CINOUT(I)
ENDDO
END

```

図3-3-13(2)

```

PROGRAM MAIN
INCLUDE 'mpif.h'
EXTERNAL PARA_SUMC16
COMPLEX*16 C, CC
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)

```

【例4】合計と最大値を1回の通信で行う

合計と最大値の両方を求めるループを並列化した場合、ループ終了後に、各プロセスが求めた合計の合計と最大値の最大値を求めるため、MPI\_REDUCEを2回コールする必要があります。ところが通信を1回行うたびに立ち上がり時間がかかります(4-1-2節参照)。ユーザーが「合計の合計と最大値の最大値を同時に求める演算」を定義すれば、通信回数は1回となり、速度が速くなる可能性があります(ただし効果があるかどうかはマシン環境に依存します)。データの動きとプログラム例を図3-3-14(1)(2)に示します。

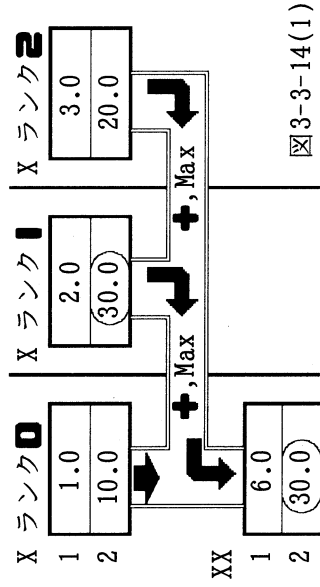


図3-3-14(1)

```

CALL MPI_OP_CREATE
& (PARA_SUMMAX,.TRUE., ISUMMAX, IERR)
X(1) = 合計を設定
X(2) = 最大値を設定
CALL MPI_REDUCE(X, XX, 1, MPI_2REAL, ISUMMAX,
& 0, MPI_COMM_WORLD, IERR)
:
SUBROUTINE PARA_SUMMAX(XIN, XINOUT, LEN, ITYPE)
REAL XIN(2,*), XINOUT(2,*)
DO I = 1, LEN
  XINOUT(1,I) = XIN(1,I) + XINOUT(1,I)
  XINOUT(2,I) = MAX(XIN(2,I), XINOUT(2,I))
ENDDO
END

```

図3-3-14(2)

```

PROGRAM MAIN
INCLUDE 'mpif.h'
EXTERNAL PARA_SUMMAX
REAL X(2), XX(2)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)

```

### 3-4 1対1通信

1対1通信サブルーチンは、2つのプロセス間での通信を行います。1対1通信サブルーチンは、本当に2つのプロセス間のみで通信を行なうために使用する場合があります。集団通信サブルーチンの機能を補う目的で、集団通信サブルーチンの代用として使用する場合があります(このときは全プロセスが1対1通信を行います)。なお、マシニング環境によっては本節の説明と動作が異なる場合があります。

#### 3-4-1 1対1通信の全体の流れ

まず1対1通信の動作の概要を図3-4-1に示します。

この例ではランク0のプロセスからランク1のプロセスにメッセージを送信しています。送信側アプリケーションはMPI\_SENDなどの送信サブルーチンをコールし、受信側アプリケーションはMPI\_RECVなどの受信サブルーチンをコールします。ただしSPMDモデルのプログラムでは、プログラム自体は1つなので、1つのプログラム内に送信と受信の両方の動作を記述します。

プログラムで使用する変数または配列である送信バッファ、受信バッファとは別に、MPIの通信サブルーチンは内部的にシステムバッファを使用します。システムバッファもメモリー上に作成されますが、アプリケーションプログラムから見ることではできません。

送信側アプリケーションプログラムは、まず①で送信バッファにデータ(メッセージ)をセットし、②でMPIの送信サブルーチン(例えばMPI\_SEND)をコールします。MPI\_SENDは③でメッセージをいったんシステムバッファにコピーし、④で宛先プロセスに送信します。

受信側アプリケーションプログラムは、まず①でMPIの受信サブルーチン(例えばMPI\_RECV)をコールします。MPI\_RECVは②で送信元プロセスからのメッセージをいったんシステムバッファに受信し、③で受信バッファにコピーします。最後に④で、受信側アプリケーションプログラムは受信バッファを参照します。

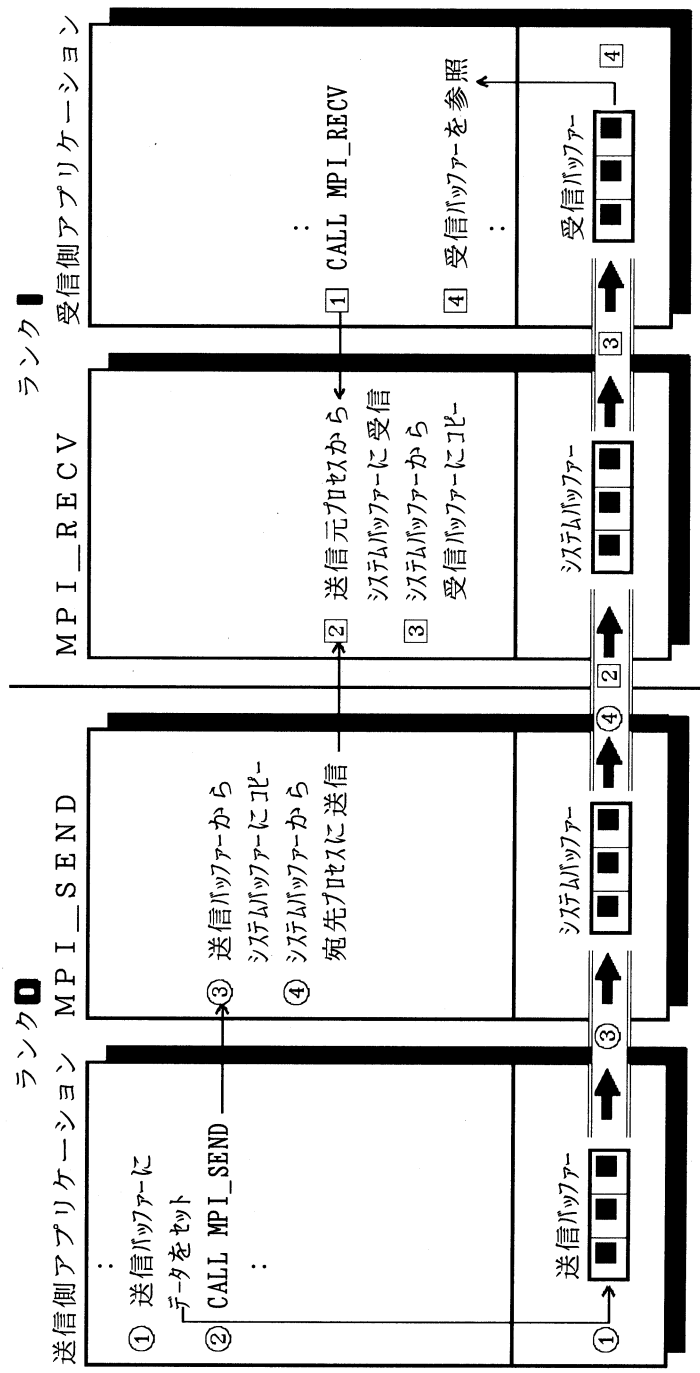


図3-4-1

## ■ 1対1通信サブルーチンの分類

図3-4-1で説明したように、1対1通信サブルーチンには送信サブルーチンと受信サブルーチンがあり、それぞれとは別に、1対1通信サブルーチンは、ブロッキング通信と非ブロッキング通信の2種類に分類することができます。ここでいうブロッキングとは、塊という意味ではなく、野球でキャッチャーが三塁からの走者をブロックする、すなわち阻止するという意味です。

以上を合わせると、MPIで提供されている主要な1対1通信サブルーチンは右のように分類されます。

実は、MPIではこれ以外にも多くの1対1通信サブルーチンが提供されています。しかし、よほど特殊な用途以外では、右のサブルーチンだけで十分だと思われまので、本書では説明を省略します。

	送信	受信
ブロッキング通信	MPI_SEND	MPI_RECV
非ブロッキング通信	MPI_ISEND + MPI_WAIT	MPI_IRECV + MPI_WAIT

## ■ ブロッキング送信(MPI\_SEND)

まず、ブロッキング送信サブルーチンMPI\_SENDについて説明します。ここで、付録のMPI\_SENDの説明を一読して下さい...

付録と同じ図を図3-4-2に示します。④でアプリケーションプログラムは『待ち』の状態に入りますが、もし『待ち』に入らないとどうなるか考えてみましょう。

通常、③のMPI\_SENDが一度だけコールされることはあまりなく、以下のようにMPI\_SENDが何度もコールされます。もし④で『待ち』の状態に入らないとすると、アプリケーションプログラムはどんどん進み、ルーブが一周して再び②を実行します。ところがこの時点で、MPI\_SENDが前回の①を完了している保証はありません。もし完了していないうちに2回目の②を実行すると、送信バッファ内にある前回の(仕掛かり中の)データを壊してしまい、誤ったメッセージが送信されてしまいます。

ブロッキング送信サブルーチンMPI\_SENDでは、このような問題が発生しないように、③を実行した後、④でアプリケーションプログラムを『待ち』の状態にし、先に進むのをブロック(阻止)します。そして、①で全てのメッセージを送信バッファbufからシステムバッファにコピーしたら、アプリケーションプログラムの『待ち』を解除します。

なお、付録にも書きましたが、『待ち』が解除されたということは、送信したメッセージが宛先プロセスに到着したことを意味するわけではありません。宛先プロセスがまだMPI\_RECV(またはMPI\_IRECV)をコールしていない場合があります。次のページで説明するMPI\_ISENDの場合も同様です。

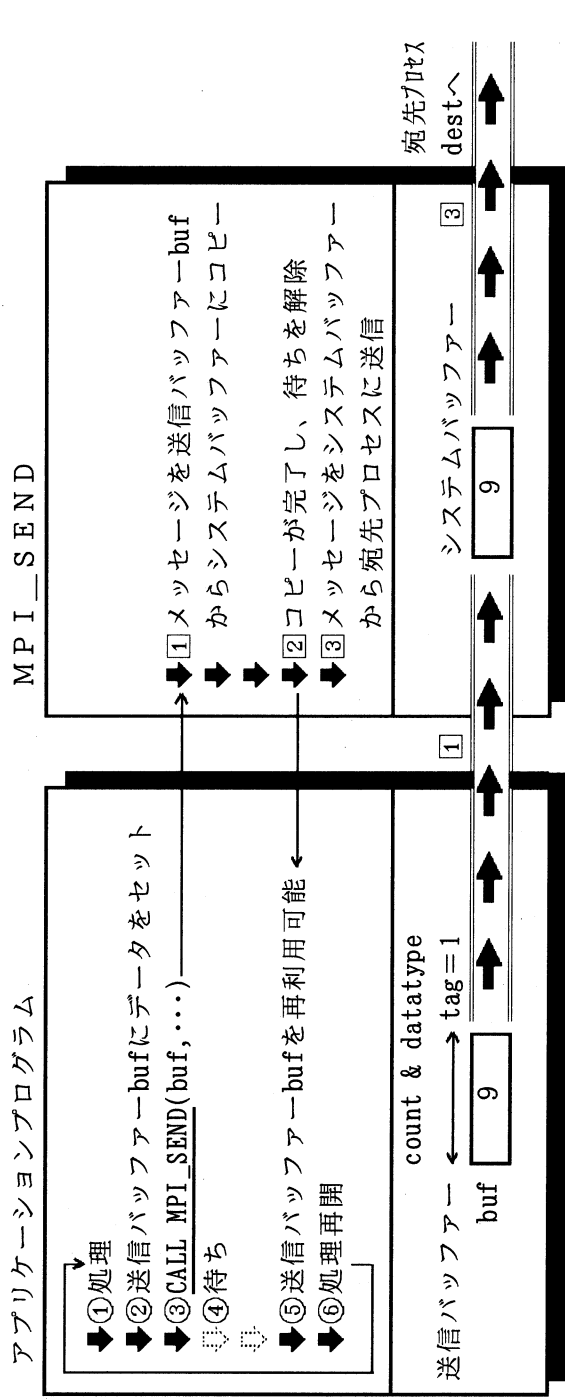


図3-4-2

■ 非ブロッキング送信(MPI\_ISEND)

次に、非ブロッキング送信サブルーチンMPI\_ISENDとMPI\_WAITを説明します。付録のMPI\_ISENDとMPI\_WAITの説明を一読して下さい。。。

付録と同じ図を図3-4-3に示します。アプリケーションプログラムは③でサブルーチンMPI\_ISENDに対してメッセージを『送れ!』と指示した後、④に示すようにほとんど先に進みます。しかし、これではループが一周して再び②を実行したときに、MPI\_SENDの説明で述べた『送信バッファ一内にある前回の(仕掛けり中の)データを壊してしまう』という問題が発生してしまいます。

これを防ぐため、2回目の②を実行するよりも前のどこかで、⑤に示すようにサブルーチンMPI\_WAITをコールします。するとアプリケーションプログラムは⑥で『待ち』に入り、MPI\_ISENDが②を実行するまで『待ち』の状態が続きます。つまりサブルーチンMPI\_WAITはストッパーの役割を果たします。

もう少し正確に述べると、③をコールすると、『送れ!』と指示したメッセージに対して識別子(ID)がMPI\_ISENDによって付けられ、その値が③の引数requestに戻ります。この値を⑤の最初の引数で指定すると、そのメッセージに対して②が実行されるまで⑥で『待ち』の状態になります。つまり引数requestは③と⑤を対応付ける役割を持ちます。

付録にも書きましたが、MPI\_ISENDを指定してMPI\_WAITを指定し忘れた場合、タイミングによる(再現性のない)エラーが発生します。またMPI\_ISENDとMPI\_WAITの引数requestがスベルミスで異なっている場合も、当然ながらMPI\_WAITは働かず、タイミングによる(再現性のない)エラーが発生しますので、スベルミスには十分注意して下さい(一方のスベルをカット・アンド・ペーストで他方にコピーするのが安全です)。

MPI\_ISENDをコールした直後にMPI\_WAITをコールした場合、MPI\_SENDをコールしたのと機能的には同じこととなります。それでは一体何のためにMPI\_ISENDがあるのでしょうか? 以下の2つの理由がありますが、本書では(1)の理由で使用します。

- (1) MPI\_SENDを使用した場合、書き方が悪いと、デッドロックという現象が発生してプログラムが止まってしまう可能性があります。これに対し『MPI\_ISEND+MPI\_WAIT』という2つの機能に分かれていると、デッドロックを起こさないプログラムを書くことができます(3-4-4節参照)。
- (2) 計算を行うCPUと独立に動作する通信専用のCPUが付いているマシン環境の場合、図3-4-3の⑤を④のずっと先(ただし2回目の②よりは前)に置くと、③の通信と④の計算がオーバーラップして実行し、通信時間が隠蔽されるため、通信のオーバーヘッドを減らす事ができます(オーバーヘッド削減)。

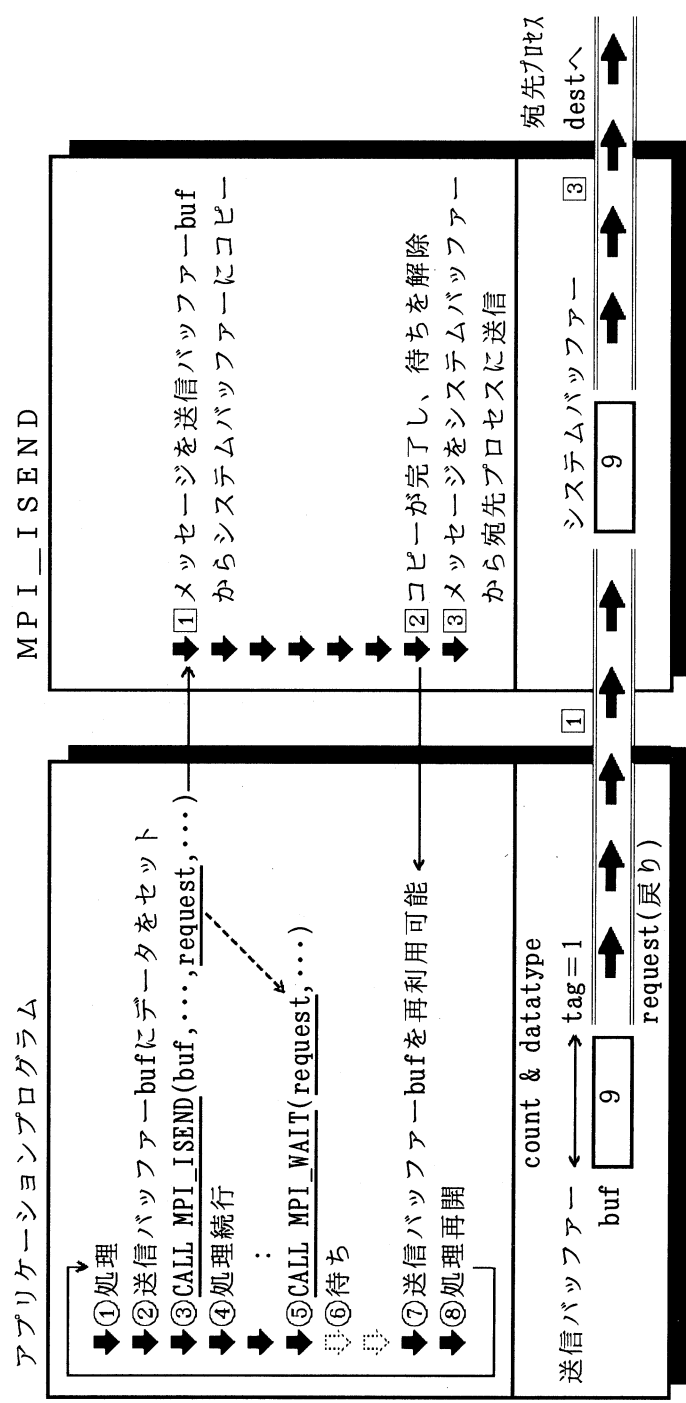


図3-4-3



■ ブロック受信(MPI\_RECV)

次に、ブロック受信サブルーチンMPI\_RECVを説明します。付録のMPI\_RECVの説明を一読して下さい  
 . . . .

付録と同じ図を図3-4-4に示します。③でアプリケーションプログラムは『待ち』の状態に入りますが、もし『待ち』に入らないとどうなるか考えてみましょう。  
 メッセージを受信するということは、受信した後にメッセージを参照するという目的があるはずで、以下の図では⑥がそれに該当します。もし③で『待ち』の状態に入らないとすると、アプリケーションプログラムはどんどん進み、⑥で受信バッファbufを参照します。ところがこの時点で、MPI\_RECVが④を完了している保証はありません。もし完了していないうちに⑥を実行すると、誤ったメッセージを参照してしまいます。

ブロック受信サブルーチンMPI\_RECVでは、このような問題が発生しないように、②を実行した後、③でアプリケーションプログラムを『待ち』の状態にし、先に進むのをブロック(阻止)します。そして、②で全てのメッセージをシステムバッファから受信バッファbufにコピーしたら、アプリケーションプログラムの受信バッファbufを参照することが可能になるので、③でアプリケーションプログラムの『待ち』を解除します。

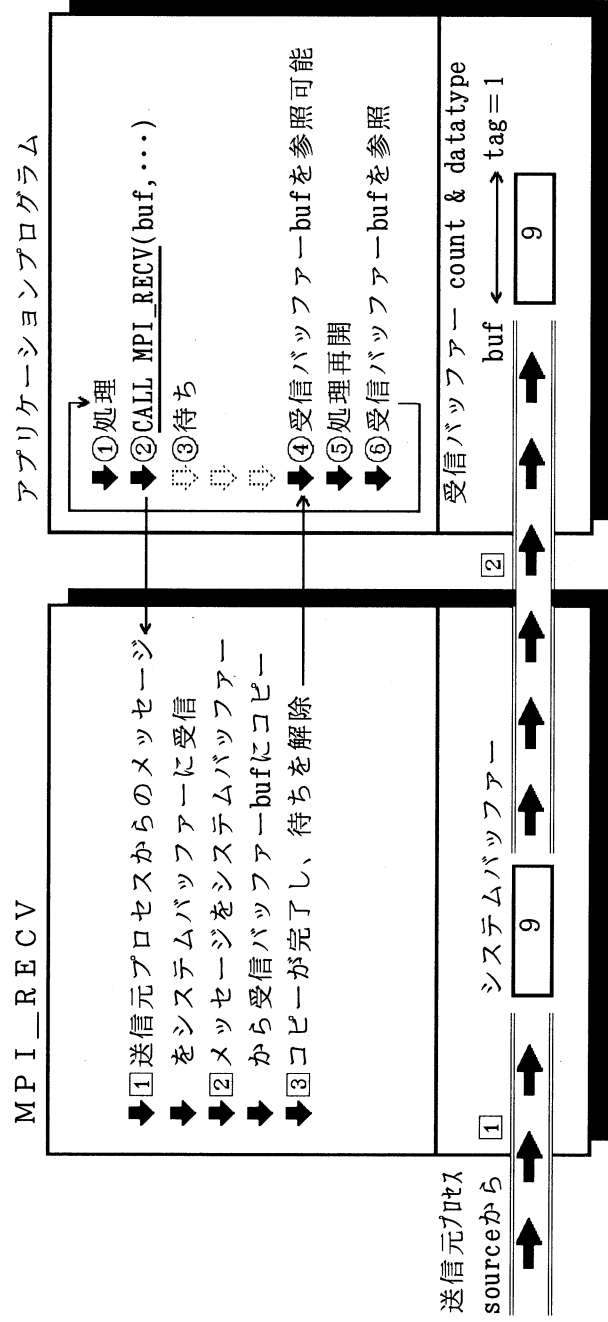


図3-4-4

## ■ 非ブロッキング受信(MPI\_I\_RECV)

最後に、非ブロッキング受信サブルーチンMPI\_I\_RECVとMPI\_WAITを説明します。付録のMPI\_I\_RECVとMPI\_WAITの説明を一読して下さい。。。

付録と同じ図を図3-4-5に示します。アプリケーションプログラムは②でサブルーチンMPI\_I\_RECVに対してメッセージを『受け取れ!』と指示した後、③に示すようにどんどん先に進みます。しかし、これでは⑥を実行したときに、MPI\_RECVの説明で述べた『受信バッファbufを参照したときに正しい正しいメッセージが入っていない』という問題が発生してしまいます。

これを防ぐため、⑧を実行するよりも前のどこかで、④に示すようにサブルーチンMPI\_WAITをコールします。するとアプリケーションプログラムは⑤で『待ち』に入り、MPI\_RECVが⑥を実行するまで『待ち』の状態が続きます。つまりサブルーチンMPI\_WAITはストッパーの役割を果たします。

もう少し正確に述べると、②をコールすると、『受け取れ!』と指示したメッセージに対して識別子(ID)がMPI\_I\_RECVによって付けられ、その値が②のindexrequestに戻ります。この値を④の最初のindexで指定すると、そのメッセージに対して③が実行されるまで⑤で『待ち』の状態になります。つまりindexrequestは②と④を対応付ける役割を持ちます。

付録にも書きましたが、MPI\_I\_RECVを指定してMPI\_WAITを指定し忘れた場合、タイミングによる(再現性のない)エラーが発生します。またMPI\_I\_RECVとMPI\_WAITのindexrequestがスペルミスで異なっている場合も、当然ながらMPI\_WAITは働かず、タイミングによる(再現性のない)エラーが発生しますので、スペルミスには十分注意して下さい(一方のスペルをカット・アンド・ペーストで他方にコピーするのが安全です)。

MPI\_I\_RECVをコールした直後にMPI\_WAITをコールした場合、MPI\_RECVをコールしたのと機能的には同じになります。MPI\_RECVでなく『MPI\_I\_RECV+MPI\_WAIT』を使用する理由は、MPI\_I\_RECVで説明したのと同様に、2つの機能に分かれていることによってデッドロックを起ささないプログラムを書くことができるため(3-4-4節参照)と、計算と通信をオーバーラップさせて通信時間を隠蔽し、通信のオーバーヘッドを減らすため(ただし計算を行うCPUと独立に動作する通信専用のCPUが付いているマシン環境の場合)です。

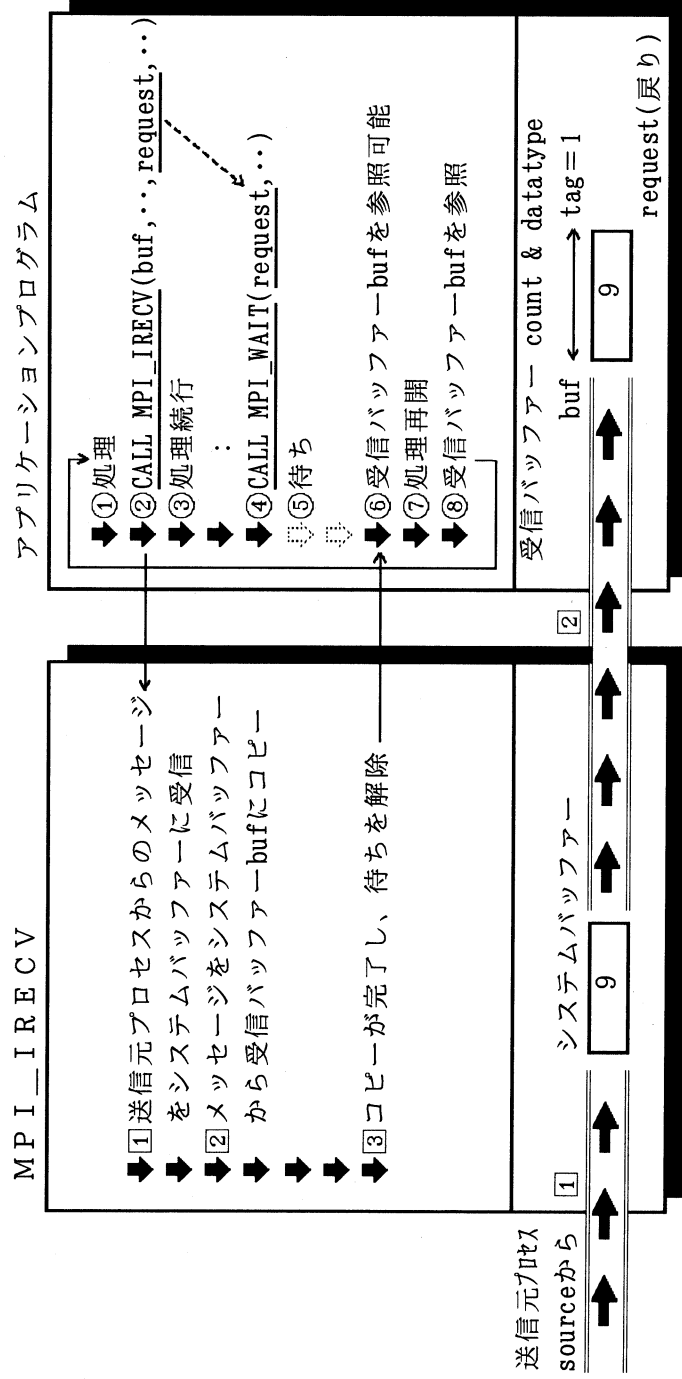


図3-4-5

### 3-4-3 1方向のみの通信

本節からは、1対1通信サブルーチンを使用して実際に通信を行う際の考慮点について説明します。まず図3-4-6(1)のように、ランク0のプロセスの送信バッファISENDからランク1のプロセスの受信バッファIRECVに(一方向に)送信する場合について説明します。

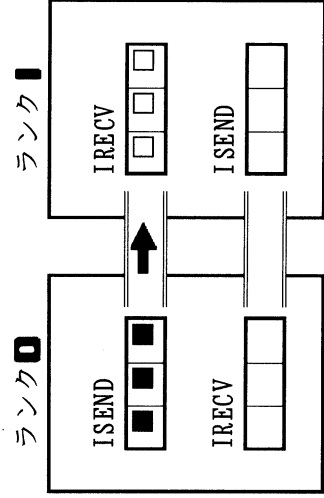


図3-4-6(1)

SPMDモデルではプログラムは1つなので、1つのプログラム内に、送信と受信の両方の動作を記述します。MPI\_SENDで送信してもMPI\_ISENDで送信しても、同じメッセージが宛先プロセスに届きます。同じメッセージなので、MPI\_RECVでもMPI\_IRECVでも受信することができます。

従って、MPI\_SENDとMPI\_ISENDのどちらで送信した場合でも、MPI\_RECVとMPI\_IRECVのどちらで受信しても構いません。すなわち送受信には、図3-4-6(2)に示すように全部で4通りの組み合わせがあります。

	受信がMPI_RECV	受信がMPI_IRECV
送信がMPI_SEND	<pre> IF (MYRANK == 0) THEN   CALL MPI_SEND(ISEND, ~) ELSEIF (MYRANK == 1) THEN   CALL MPI_RECV(IRECV, ~) ENDIF : </pre>	<pre> IF (MYRANK == 0) THEN   CALL MPI_SEND(ISEND, ~) ELSEIF (MYRANK == 1) THEN   CALL MPI_IRECV(IRECV, ~, IREQ, ~)   CALL MPI_WAIT(IREQ, ~) ENDIF : </pre>
送信がMPI_ISEND	<pre> IF (MYRANK == 0) THEN   CALL MPI_ISEND(ISEND, ~, IREQ, ~)   CALL MPI_WAIT(IREQ, ~) ELSEIF (MYRANK == 1) THEN   CALL MPI_RECV(IRECV, ~) ENDIF : </pre>	<pre> IF (MYRANK == 0) THEN   CALL MPI_ISEND(ISEND, ~, IREQ, ~) ELSEIF (MYRANK == 1) THEN   CALL MPI_IRECV(IRECV, ~, IREQ, ~) ENDIF CALL MPI_WAIT(IREQ, ~) : </pre>

図3-4-6(2)

### 3-4-4 双方向通信

#### ■ 双方向通信とデッドロック

次に、図3-4-7に示すように、ランク0のプロセスからランク1のプロセスに送信し、ランク1からランク0にも送信するという、双方向の通信について説明します。ここでは双方向とは、たまたまプログラム内のほぼ同じ場所で、ランク0からランク1へのメッセージの送信と、ランク1からランク0へのメッセージの送信が必要になったという意味であると考えて下さい。

図3-4-8に示すように、2人が相手の家に引越すとします。その時2人とも、相手の家が空いたら荷物を運ぼうと考えたとすると、2人とも動けなくなってしまいます。このように、お互いに相手の動作待ちになり、こう着状態に陥ることをデッドロックと呼びます。

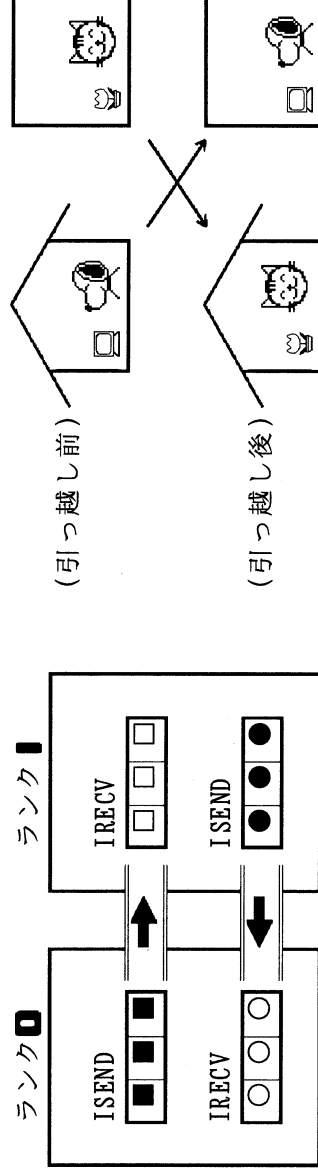


図3-4-8

双方向に通信を行う場合、通信順序がまずいとデッドロックが発生し、上記の例と同様に、各プロセスは相手プロセスの動作待ちになり、両方のプロセスが動けなくなってしまいます。デッドロックが発生すると、プログラムはその箇所で永久に止まったままになるため、並列ジョブをキャンセルするしかありません。これは、前述の再現性のないエラーと並んで、並列ジョブに特有の現象です。

図3-4-7のように双方向に通信を行う場合、各プロセスは相手のプロセスに対して送信と受信の両方を行います。各プロセスの送信と受信の実行順序によって、図3-4-9に示す3つのパターンが考えられます。プロキキングと非プロキキングのどちらかで通信するかまで考慮すると、さらにバリエーションが増えます。なお、図中の「MPI\_(I)SEND」は、MPI\_SENDとMPI\_ISENDのいずれかであることを示します。また、この図ではMPI\_WAITは省略しました。

このように、双方向通信には多くのバリエーションが考えられます。それでは、図3-4-9のそれぞれのパターンについてデッドロックの可能性を検討してみましょう。

パターン1

```

:
:
IF (MYRANK == 0) THEN
  CALL MPI_(I)SEND(~)
  CALL MPI_(I)RECV(~)
ELSEIF (MYRANK == 1) THEN
  CALL MPI_(I)SEND(~)
  CALL MPI_(I)RECV(~)
ENDIF
:
:

```

ランク0、1のプロセスがともに送信、受信の順にコール。

パターン2

```

:
:
IF (MYRANK == 0) THEN
  CALL MPI_(I)RECV(~)
  CALL MPI_(I)SEND(~)
ELSEIF (MYRANK == 1) THEN
  CALL MPI_(I)RECV(~)
  CALL MPI_(I)SEND(~)
ENDIF
:
:

```

ランク0、1のプロセスがともに受信、送信の順にコール。

パターン3

```

:
:
IF (MYRANK == 0) THEN
  CALL MPI_(I)SEND(~)
  CALL MPI_(I)RECV(~)
ELSEIF (MYRANK == 1) THEN
  CALL MPI_(I)RECV(~)
  CALL MPI_(I)SEND(~)
ENDIF
:
:

```

ランク0のプロセスは送信、受信、ランク1のプロセスは受信、送信の順にコール  
(ランク0と1は逆でも可)。

図3-4-9

■ パターン1

最初に結論を述べます。図3-4-10(1)のように、両プロセスがともにブロッキングの送信MPI\_SEND、受信の順に通信を行った場合、デッドロックが発生する可能性があります。図3-4-10(2)の②のように非ブロッキング送信MPI\_ISENDの直後にMPI\_WAITをコールした場合も、①と機能的に同じなのでデッドロックが発生します。図3-4-10(3)のように、③に対するMPI\_WAITを④の後ろの⑤に置いた場合は、デッドロックにはなりません(理由は後述します)。

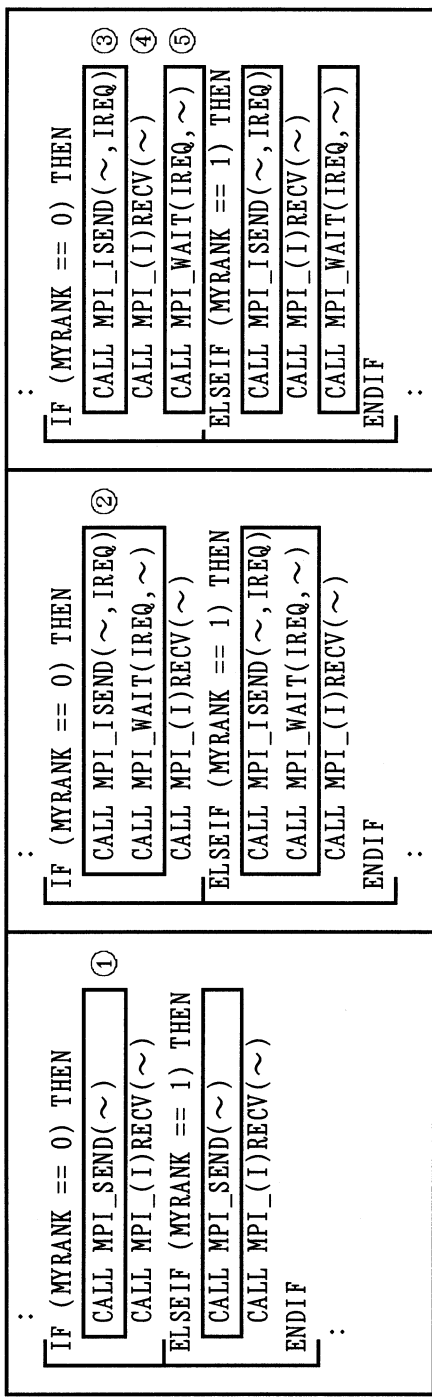


図3-4-10(1) ✖

図3-4-10(2) ✖

図3-4-10(3) ○

デッドロックが発生する理由を説明します(マシン環境によってはメカニズムが異なるため、発生しない場合もあります)。図3-4-11で、送信メッセージの長さ「4」が、システムバッファアの大きさ「2」より長いです。この場合、ランク0のMPI\_SENDは、まず■をランク1に送信し、次に残りの□を送信します。■はMPI\_SENDのプログラムは①をコールすると、②で待ちに入ります。コールされたMPI\_SENDは、(1)です。システムバッファアにコピーします。(2)の時点でランク1のプログラムが④のMPI\_I(IRECV)をすでにコールしている、(3)で■を送信できますが、まだコールしていない場合は④のMPI\_I(IRECV)を待ちます。

一方ランク1のプログラムとMPI\_SENDもランク0と同じ状態になっていて、それぞれ②と②で止まっています。このため④のMPI\_I(IRECV)がコールされず、前述のランク0のMPI\_SENDは(2)で止まったまま先に進みません。すなわちランク0、1のプログラムとMPI\_SENDは、全員下線部分で止まり、先に進めなくなります。このデッドロックは、④と④のMPI\_I(IRECV)を動かすことができれば回避することができます。MPI\_SENDはMPI\_ISENDとMPI\_WAITという2つのルーチンに分けることができます。そこで図3-4-10(3)のようにすれば、MPI\_ISENDの他にMPI\_I(IRECV)も動くので、デッドロックを回避することができます。

なおこのデッドロックは「誤り」ですから作成しないようにして下さい。例えば、並列化したプログラムがテスト段階では小規模なデータを使用したため、たまたまデッドロックにならず、本番で大規模なデータを使用したらデッドロックになったり、あるいはマシン環境を変えたらデッドロックになる恐れがあります。

MPI\_SEND 【ランク0】 MPI\_ISEND

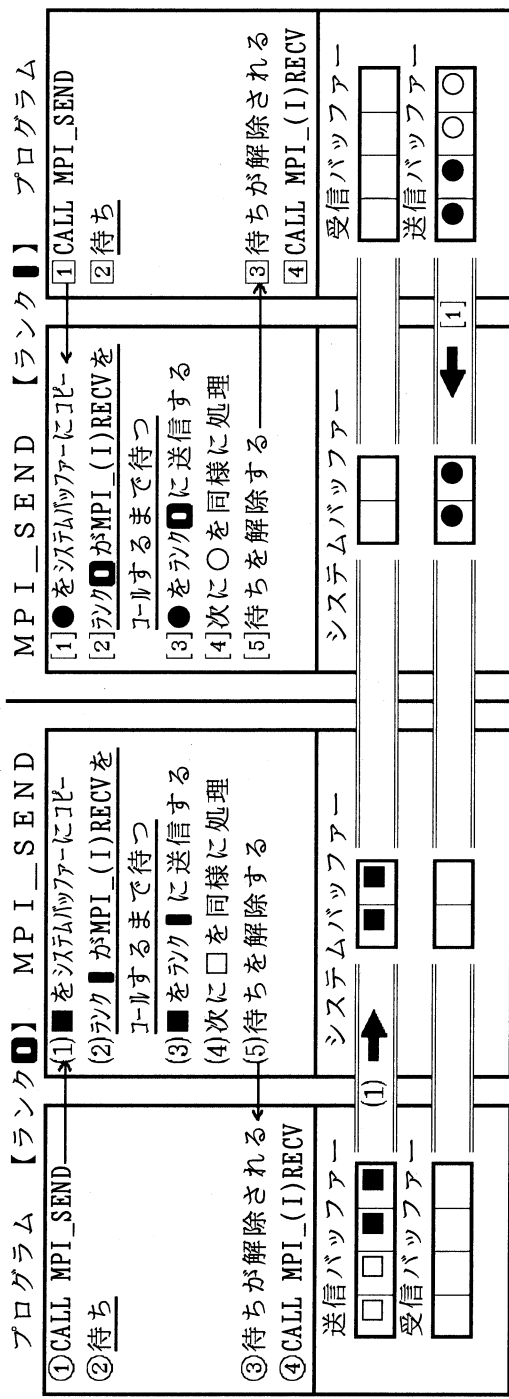


図3-4-11

最初に結論を述べます。図3-4-12(1)のように、両プロセスがともにブロッキングの受信MPI\_RECV、送信の順に通信を行った場合、必ずデッドロックが発生します。図3-4-12(2)の②のように非ブロッキング受信MPI\_RECVの直後にMPI\_WAITをコールした場合も、①と機能的に同じなのでデッドロックが発生します。図3-4-12(3)のように、③に対するMPI\_WAITを④の後ろの⑤に置いた場合は、デッドロックにはなりません(理由は後述します)。

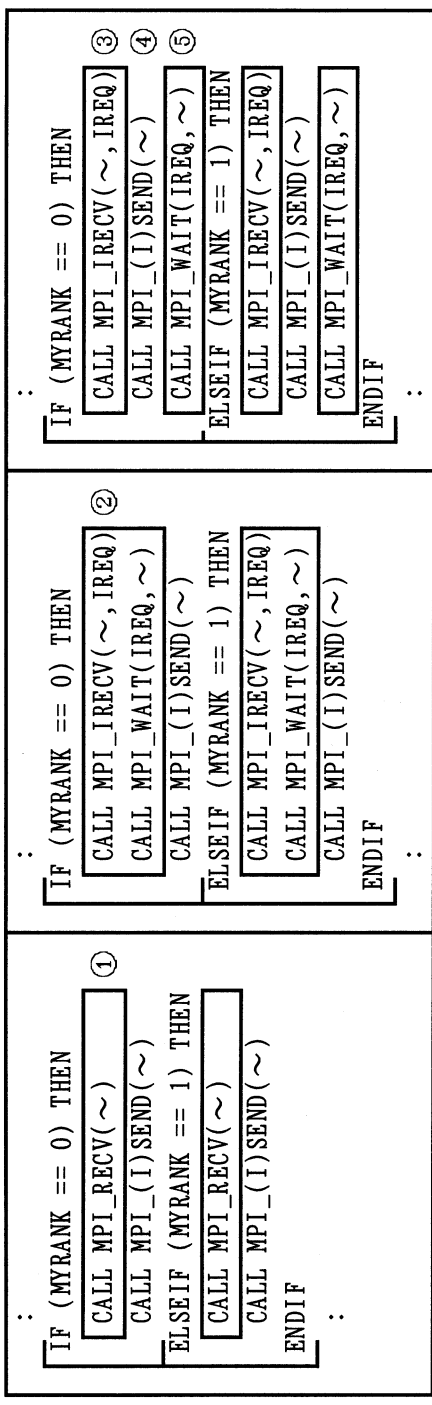


図3-4-12(1) ✕ 図3-4-12(2) ✕ 図3-4-12(3) ○

このデッドロックは、前述のパターン1のデッドロックに比べると、発生メカニズムは簡単です。

図3-4-13に示すように、ランク0のプログラムは①をコールすると、②で待ちに入ります。コールされたMPI\_RECVは、(1)の時点でランク1のプログラムが④のMPI\_RECVをすでにコールしていれば、(2)で受信できず、まだコールしていない場合はコールするまで(1)で待ちます。

一方ランク1のプログラムとMPI\_RECVもランク0と同じ状態になっていて、それぞれ②と①で止まっています。このため④のMPI\_RECVがコールされず、前述のランク0のMPI\_RECVは(1)で止まったまま先に進みません。すなわちランク0、1のプログラムとMPI\_RECVは、全員下線部分で止まり、先に進めなくなります。このデッドロックは、④と④のMPI\_RECVを動かすことができれば回避することが出来ます。MPI\_RECVはMPI\_RECVとMPI\_WAITという2つのルーチンに分けることができます。そこで図3-4-12(3)のようにすれば、MPI\_RECVの他にMPI\_RECVも動くので、デッドロックを回避することができます。

パターン1と違い、このデッドロックは送信メッセージが短くても必ず発生します。実際にこのデッドロックが発生するプログラム例を、4-6-2節の「境界条件の処理」に示します。

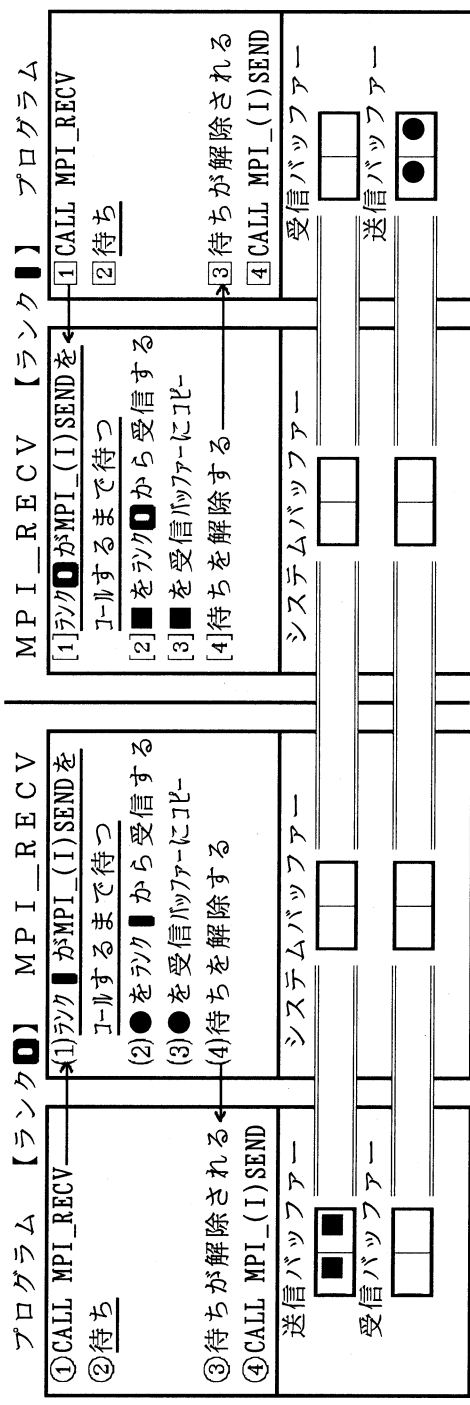


図3-4-13

■ パターン3

前述の図3-4-7の双方向通信は、図3-4-6(1)の一方向通信を、向きを変えて2回行ったとみなすことができ、プログラムで書くと図3-4-14(1)のようになります(図でMPI\_WAITは省略しています)。まず①でランク0のプロセスからランク1のプロセスに一方向通信を行い、次に②でランク1のプロセスからランク0のプロセスに一方向通信を行います。

①と②のIF文を整理すると図3-4-14(2)になります。これをよく見ると、図3-4-9のパターン3と同じです。そして図3-4-6(2)で述べたように、一方向通信はプロッキングと非プロッキングの送受信のいずれの組み合わせでも可能なので、図3-4-14(1)(2)もいずれの組み合わせでも可能です。

まとめると、パターン3は、プロッキングと非プロッキングのいずれの通信ルーチンの組み合わせでもデッドロックにはなりません。

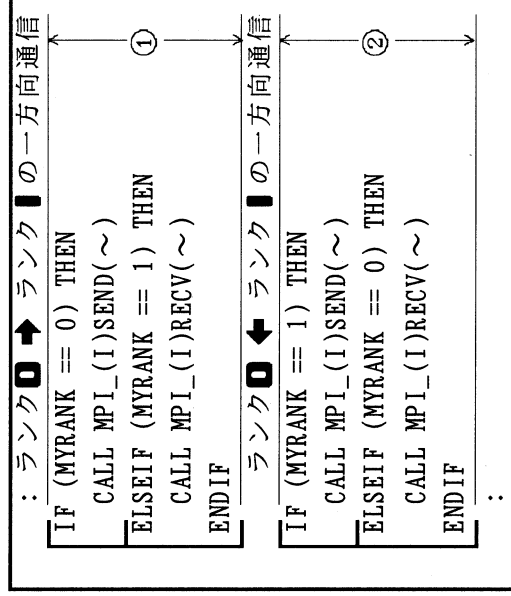


図3-4-14(1)  
1方向通信を  
2回行う

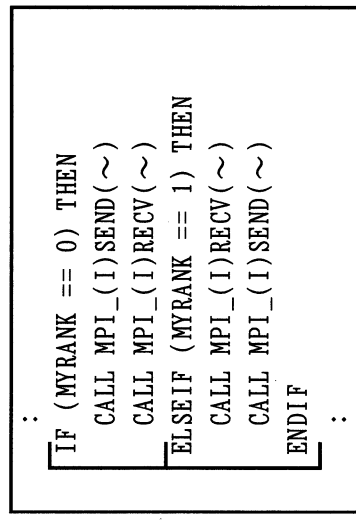


図3-4-14(2) パターン3

■ 全ての通信を非プロッキング通信にした場合

双方向通信の場合、パターン1とパターン2で説明したように、送信ルーチンや受信ルーチンが動けなくなるためにデッドロックが発生することが分かりました。従って、図3-4-15のように送信と受信ルーチンを全て非プロッキング通信にし、全ての送受信ルーチンをいっせいに動かすようにすれば、デッドロックになることはありません。

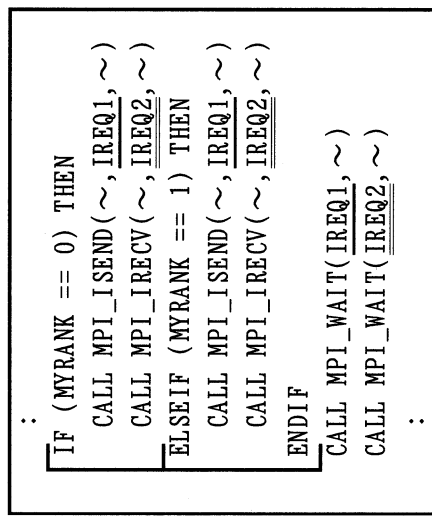


図3-4-15 ○

■ 1対1プロッキング通信ルーチンMPI\_SENDRCV

1対1通信ルーチンMPI\_SENDRCVは送信と受信が合体したルーチンで、図3-4-16(1)(2)のような通信を行うことができます(詳細は4-6-2節と付録参照)。このルーチンはプロッキング通信ルーチンですが、デッドロックにはなりません。4-6-2節、5-1節で説明する、シフトと呼ばれる通信パターンに特に有効です。分類上は1対1通信ルーチンですが、集団通信ルーチンの1つと考えることもできます。

MPI\_SEND, MPI\_RECV, MPI\_ISEND, MPI\_IRECV, MPI\_SENDRCVのどれが最も速いかはマシン環境に依存します。

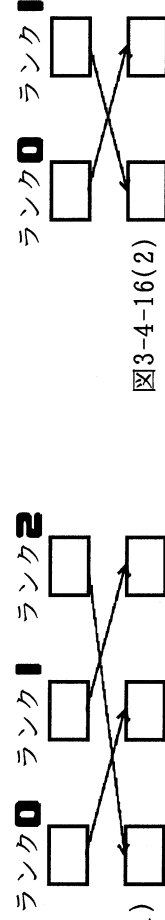


図3-4-16(1)

図3-4-16(2)

### 3-4-5 1対1通信に関する補足

いくつか補足します。

- 1対1通信ルーチンでは、受信バッファで指定するバイト数(データ個数×データ型のバイト数)は、送信バッファで指定するバイト数と等しいか、または大きくても構いません。
- 3-3-4節でも書きましたが、1対1通信および集団通信ルーチンで一度に送信できる最大のバイト数はマシンの環境によって異なり、無制限の場合と、制限のある場合があります。制限のある場合に、制限を越えたバイト数を送ってもエラーにならず、不完全なメッセージが送られてしまうマシン環境もありますので、お使いのマシンのマニュアルを確認して下さい。

● ランク0のプロセスからランク1のプロセスに1対1送信(MPI\_SENDなど)を2回連続して行った場合、メッセージはコールした順番にランク1のプロセスに到着します。言いかえると、あるメッセージが途中で前のメッセージを追い越すことはないのです。プログラムではメッセージの順序を気にする必要はありません。

● MPIで並列化したプログラムを1プロセスで実行した場合、通常は正常に動作します。しかし例えば図3-4-17のように、1対1通信ルーチンを含む並列プログラムを1プロセスで実行した場合、そのプロセスのランクは0なので①が真となり②が実行されます。しかし③は偽となり④が実行されないで、以後の計算がおかしくなってしまいます。このような問題が発生する可能性があるので、1プロセスで実行するときは、並列版ではなく単体版のプログラムを使用することを強くお勧めします(4-6-2節の「境界条件の処理」参照)。

```

:
IF (MYRANK==0) THEN
  CALL MPI_SEND(A,1,MPI_REAL,1,~)
ELSEIF (MYRANK==1) THEN
  CALL MPI_RECV(B,1,MPI_REAL,0,~)
ENDIF
:

```

図3-4-17

● 以下のような場合、仕掛かり中のメッセージがどんどん累積していきませんが、その後の動作はマシン環境によって異なりますので(例えば受信バッファがパンクして異常終了するなど)、お使いのマシンの環境のマニュアルで確認して下さい。

-図3-4-18(1)のように、ランク1からランク0に一度に複数のメッセージを送信したものの、ランク0は別の処理中のため、まだ受信できない場合。

-図3-4-18(2)のように、複数のプロセス(本例では1~3)からランク0に一度にメッセージを送信したものの、ランク0は別の処理中のため、まだ受信できない場合。

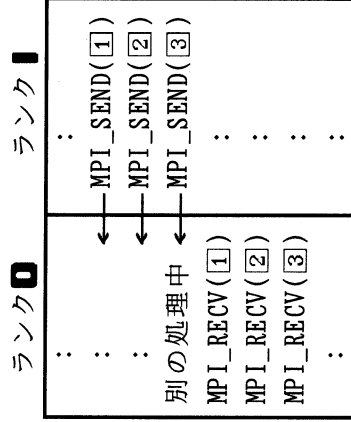


図3-4-18(1)

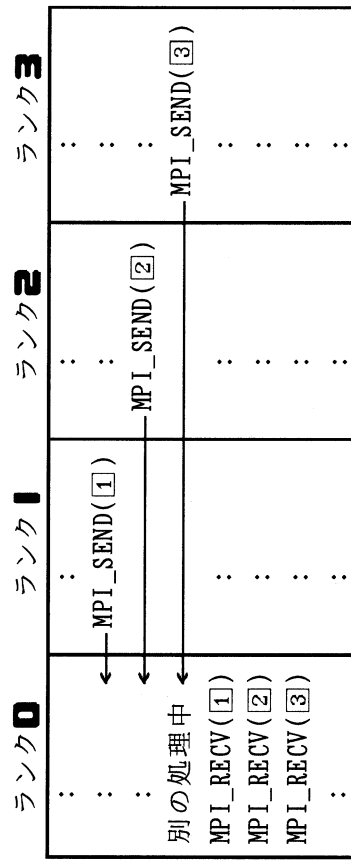


図3-4-18(2)



### 3-5 派生データ型

#### 3-5-1 派生データ型とは

図3-5-1(1)で、ランク0が持つ配列Aのうち●の部分のみを、ランク0が持つ配列Bの●の部分に送信する場合作を考えます。なお、図3-5-1(1)(2)では、前節で説明したシステムバッファは省略していただきます。注意して下さい。

今まで説明した方法では、MPIの通信ルーチンで指定する送受信バッファは、連続したデータしか扱うことができず、図3-5-1(1)のようにとびとびのデータのみを通信することはできませんでした。このためランク0のプロセスは、配列Aのうち●部分のみを一時配列TEMPにコピーしてから送信し、ランク0のプロセスは、受信したメッセージを一時配列TEMPに受信した後、配列Bにコピーしなければなりませんでした。この方法は、一時配列TEMPを新規に確保し、コピーする必要があるため面倒です。本節で紹介する派生データ型という機能を使用すると、図3-5-1(2)のように、メモリー上でとびとびのデータのみを送受信バッファを介さずに直接通信することができます。言い換えると、とびとびのデータ部分のみを送受信バッファとして指定することができます。

なお、図3-5-1(1)と図3-5-1(2)のどちらのパフォーマンスが良いかは、マシン環境に依存します。

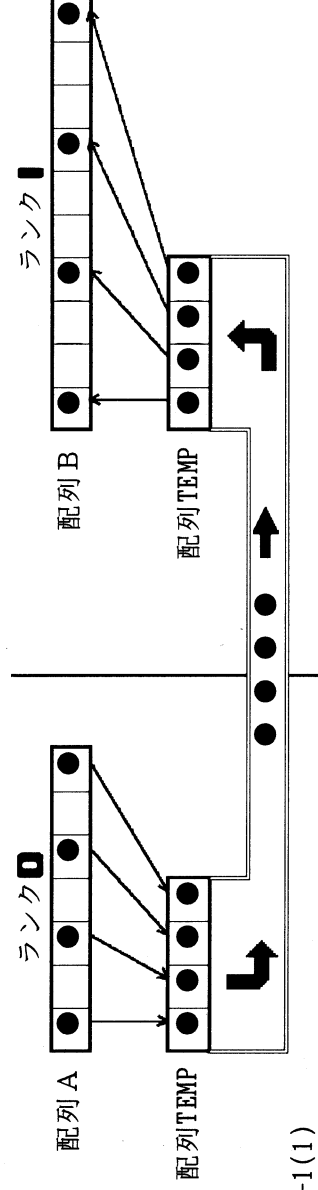


図3-5-1(1)

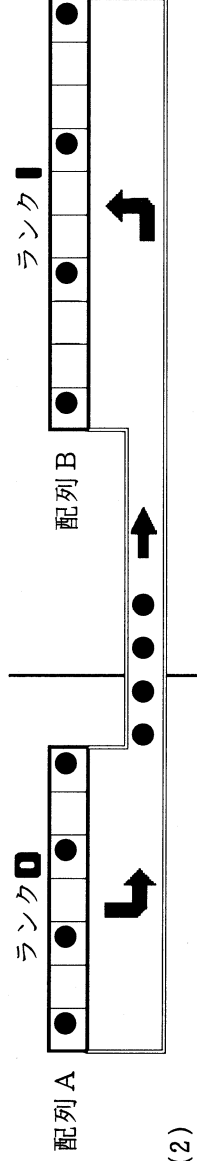


図3-5-1(2)

図3-5-2(1)に示すように、配列IBUFの中の●のデータのみを通信する場合、まず後述するMPIのルーチンを使用して、INewType1というパターンのデータ型を作成します。MPI\_INTEGERのように、MPIがあらかじめ提供しているデータ型をMPIの基本データ型と呼びます(3-3-4節参照)。これに対し、INewType1のようにユーザーが作成するデータ型を派生データ型と呼びます。

派生データ型INewType1を使用して、例えば集団通信ルーチンMPI\_BCASTで通信する場合、図3-5-2(1)の下部のように指定します。これは送信バッファが、「IBUF(4)から始まる1個のINewType1型の部分」であることを意味します。この通信を行うと、送信元のランク0のプロセスでは、配列IBUF内にとびとびに入っている6個の●のみが抽出されて送られます。宛先のランク1のプロセスでは、受け取った6個の●が配列IBUF内のとびとびの部分に直接入ります。

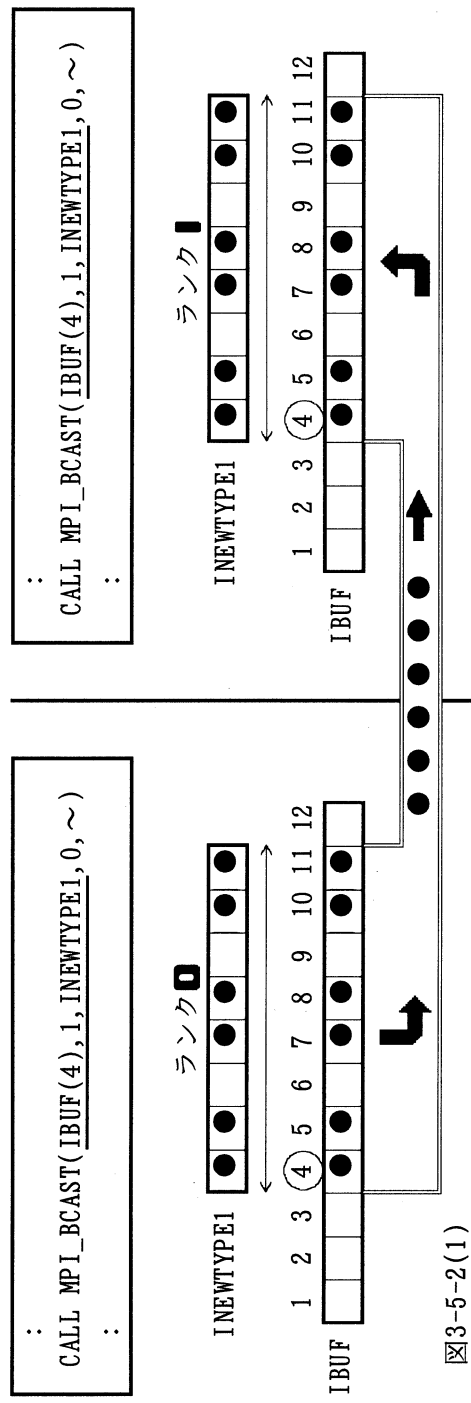


図3-5-2(1)

同様に、図3-5-2(2)のINewType2というパターンの派生データ型を作成した場合、送信バッファが、「IBUF(4)から始まる3個のINewType2型の部分」であることを意味します。

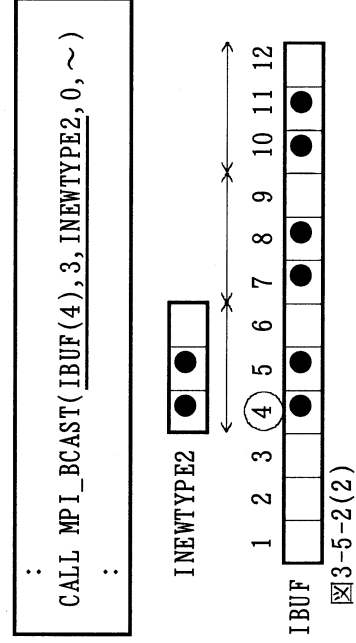


図3-5-2(2)

あまり一般的な例ではありませんが、図3-5-2(3)のように、通信ルーチンの引数で指定するデータ型がプロセスによって異なってもかまいません。ただし、送信するバイト数の合計は一致している必要があります。

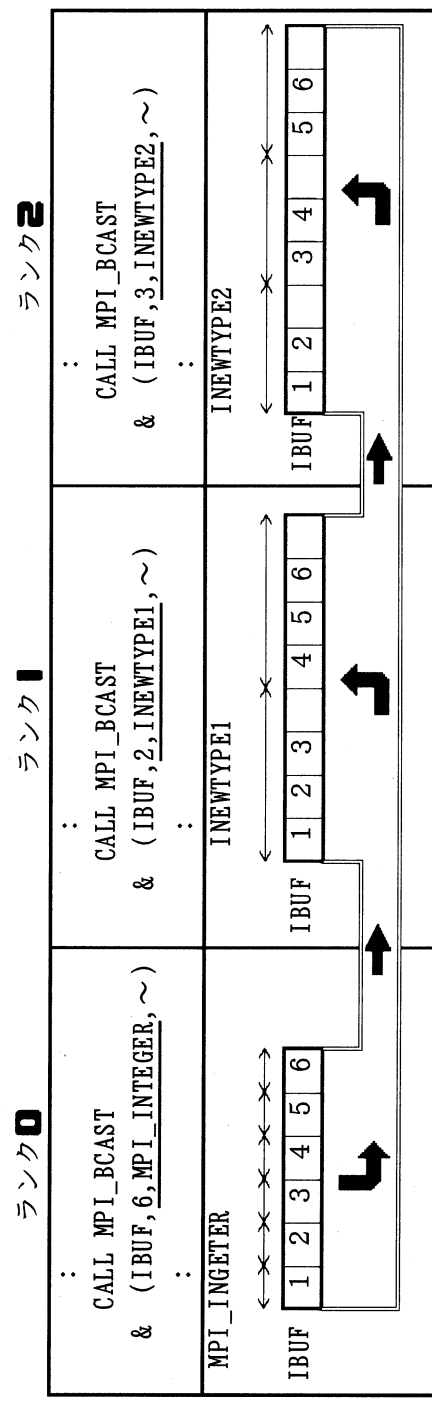


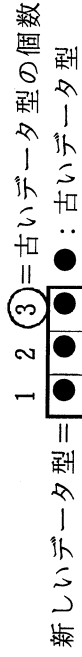
図3-5-2(3)

MPIが提供する派生データ型の主なサブルーチンを紹介します。使用方法の詳細は付録を参照して下さい。なお、付録に掲載したサブルーチンのうち、MPI\_TYPE\_SIZEの説明は、以下では省略します。

以下の説明中で、古いデータ型とは、MPIが提供している基本データ型(MPI\_INTEGERなど)、または、その時点までにユーザーが作成した派生データ型を意味します。新しいデータ型とは、古いデータ型に基づいて、以下のサブルーチンを使って新たに作成された派生データ型を意味します。

## (1) MPI\_TYPE\_CONTIGUOUS

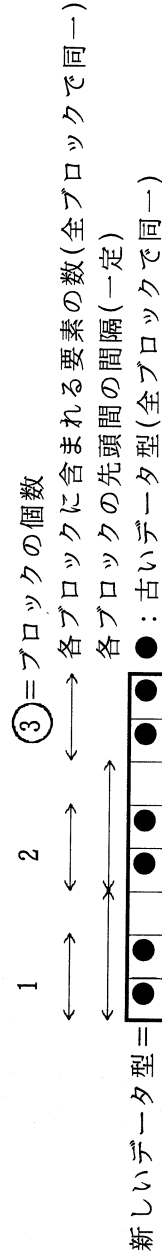
古いデータ型の要素を、指定した個数(以下では3個)だけ連結して新しいデータ型を作成します。



## (2) MPI\_TYPE\_VECTOR/MPI\_TYPE\_HVECTOR

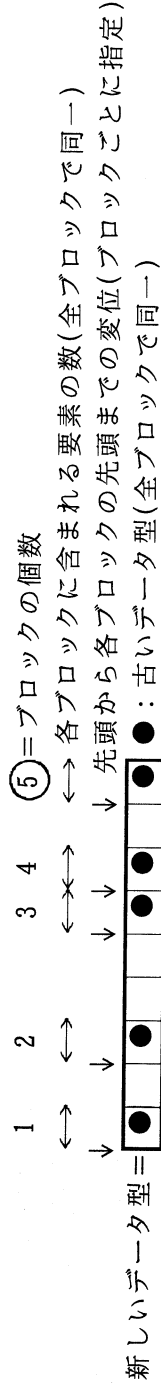
以下の(2)~(5)のサブルーチンは、複数(1つでも可)のブロックから構成される新しいデータ型を作成します。

本サブルーチンでは、各ブロックの先頭間の間隔は一定、各ブロックに含まれる古いデータ型の要素の数は同一、各ブロックに含まれる要素の古いデータ型の種類も同一です。MPI\_TYPE\_VECTORとMPI\_TYPE\_HVECTORは、ブロックの先頭間の間隔を指定する方法が異なる(要素数/バイト)以外は同じです。



## (3) MPI\_TYPE\_CREATE\_INDEXED\_BLOCK

本サブルーチンはMPI-2で提供されています。ブロックごとに、新しいデータ型の先頭から各ブロックの先頭までの変位を指定します。各ブロックに含まれる古いデータ型の要素の数は同一、各ブロックに含まれる要素の古いデータ型の種類も同一です。最初のブロックの前にブラントクがあっても構いません(後述する(7)参照)。間接アドレスを使ったリストベクトルの派生データ型を作成するときに便利です(4-6-6-1節参照)。



## (4) MPI\_TYPE\_INDEXED/MPI\_TYPE\_HINDEXED

ブロックごとに、新しいデータ型の先頭から各ブロックの先頭までの変位、各ブロックに含まれる古いデータ型の要素の数を指定します。各ブロックに含まれる要素の古いデータ型の種類は同一です。最初のブロックの前にブラントクがあっても構いません(後述する(7)参照)。MPI\_TYPE\_INDEXEDとMPI\_TYPE\_HINDEXEDは、先頭からの変位を指定する方法が異なる(要素数/バイト)以外は同じです。



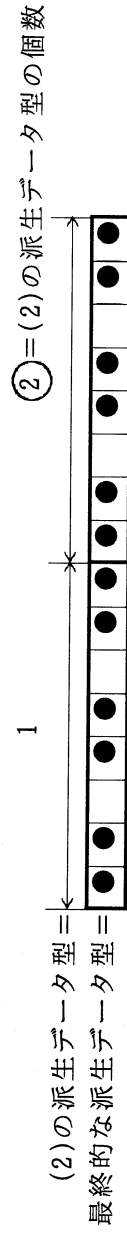
## (5) MPI\_TYPE\_STRUCT

ブロックごとに、新しいデータ型の先頭から各ブロックの先頭までの変位、各ブロックに含まれる古いデータ型の要素の数、各ブロックに含まれる要素の古いデータ型を指定します。最初のブロックの前にブランクがあっても構いません(後述する(7)参照)。本サブルーチンを使用して、構造体の派生データ型を作成することができます(3-5-4節参照)。



## (6) MPI\_TYPE\_COMMIT

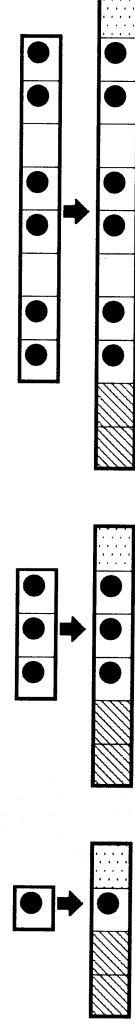
作成した派生データ型をさらに組み合わせさせて派生データ型を作成することもできます。以下の例では、まず前述の(2)の派生データ型を作成し、それを(1)で説明したMPI\_TYPE\_CONTIGUOUSで2つ連結して、最終的な派生データ型を作成しています。



最終的な派生データ型が完成した後、それをサブルーチンMPI\_TYPE\_COMMITでMPIに登録します。これが完了すると、その派生データ型を、集団通信ルーチンや1対1通信ルーチンのデータ型の引数に指定することができます。なお、作成途中の一時的な派生データ型は、MPI\_TYPE\_COMMITで登録する必要はありません。

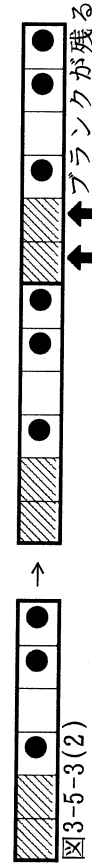
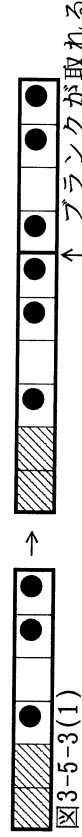
## (7) 派生データ型の前と後のブランクについて

MPIが提供する基本データ型(MPI\_INTEGERなど)や、前述の(1)(2)で作成した派生データ型の場合、以下の上図に示すように、そのままでは前後にブランクは入りません。以下の下図のように、前や後にブランク(着色部分)を入れた場合、(8)か(9)の方法で行います。



一方、前述の(3)(4)(5)で作成した派生データ型の場合、新しいデータ型の先頭から最初のブロックの先頭までの変位を0より大きな値にすると、図3-5-3(1)(左図)のように最初のブロックの前にブランクを入れることができます。ただし、この方法では派生データ型の最後のブロックの後にブランクを入れることはできません。

また、この派生データ型を使用してデータを2個以上通信する場合や、この派生データ型をMPI\_TYPE\_CONTIGUOUSで2個以上連続させた場合、図3-5-3(1)(右図)のように、2個目以降の派生データ型では前のブランクが取れてしまいます。図3-5-3(2)(右図)のようにブランクを残したい場合、(8)か(9)の方法で行います。



(8) 派生データ型の前後にブランクを入れる方法(1) (MPI\_TYPE\_STRUCTを使用)

まず、(5)で説明したMPI\_TYPE\_STRUCTのみを用いて、派生データ型の前後にブランクを入れる方法を説明します。派生データ型の最初のブロックの前にブランクを入りたい場合、図3-5-3(3)(上図)に示すように、ブランクの前に要素数1、大きさ0、データ型が「MPI\_LB」(LBはLower Boundの略)のブロックを定義します。同様に派生データ型の最後のブロックの後にブランクを入れたい場合、ブランクの後に要素数1、大きさ0、データ型が「MPI\_UB」(UBはUpper Boundの略)のブロックを定義します。

この派生データ型を使用してデータを2個以上通信する場合や、この派生データ型をMPI\_TYPE\_CONTIGUOUSで2個以上連続させた場合、図3-5-3(3)(下図)に示すように、2個目以降の派生データ型では前のブランクが残ります。

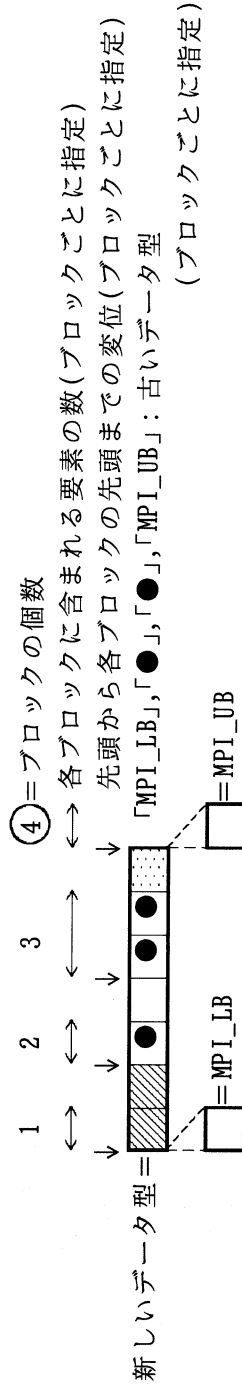
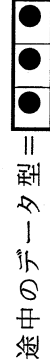


図3-5-3(3) ↑ ブランクが残る

次に、MPIが提供する基本データ型(MPI\_INTEGERなど)や、前述の(1)~(4)で作成した派生データ型の前後にブランクを入れる方法を、(1)の派生データの例で説明します。まず図3-5-3(4)(上図)に示すように、派生データ型(以後「途中のデータ型」を作成します。次に図3-5-3(4)(下図)に示すように、MPI\_TYPE\_STRUCTを使用し、上記で説明した「MPI\_LB」と「MPI\_UB」をブランクの前後に指定して「最終のデータ型」を作成します。このとき、「途中のデータ型」の要素の数は「3」ではなく「1」となります。なお、この方法で作成した派生データ型も、図3-5-3(3)と同様に、2個目以降の派生データ型では前のブランクが残ります。



③ = ブロックの個数

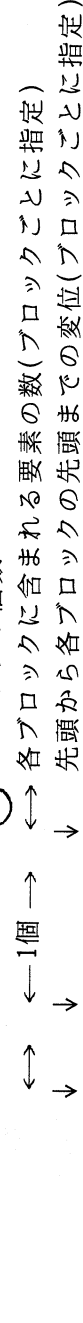


図3-5-3(4) ↑ ブランクが残る

(9) 派生データ型の前後にブランクを入れる方法(2) (MPI\_TYPE\_CREATE\_RESIZEDを使用)

MPI-2で提供されたMPI\_TYPE\_CREATE\_RESIZED(付録参照)を使って、派生データ型の前後にブランクを入れる方法を説明します。本サブルーチンは、MPIが提供する基本データ型(MPI\_INTEGERなど)や、前述の(1)~(5)で作成した派生データ型に対して「範囲」を指定します。元の大きさより大きな範囲を指定すると、図3-5-3(5)に示すように、派生データ型の後にブランクを入れます。元の大きさより大きな範囲を指定すると、この方法で派生データ型の前と後の両方にブランクを入れます。次のようにします。まず(7)で説明した方法で、(3)(4)(5)を使って図3-5-3(6)(上図)に示すように前のブランクを入れます。この派生データ型に対して、図3-5-3(6)(下図)に示すように、MPI\_TYPE\_CREATE\_RESIZEDで元の大きさと同じ大きさの範囲を指定すると、前と後の両方にブランクが入ります。2個目以降の派生データ型では前のブランクが残ります。

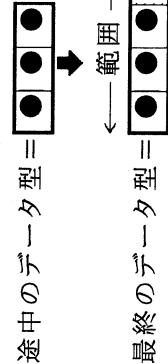


図3-5-3(5) ↑ ブランクが取れる

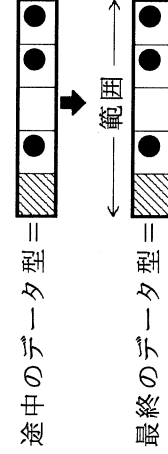


図3-5-3(6) ↑ ブランクが残る

### 3-5-3 部分配列を表す派生データ型の作成

本節と次節では、並列化したときに派生データ型が必要になる例を紹介しします。図3-5-4(1)で、ランク0が着色した部分を担当しています。Fortranでは、2次元配列はメモリー上で1次元方向に連続(C言語では逆)になっています(3-1節参照)。この例では計算領域を2次元方向に分割しており、着色した部分はメモリー上で連続なので直接送信バッファとして送信することができます。ただし、図3-5-4(2)のように周囲にノリシロがある場合、着色した部分は不連続になるので、着色した部分のみを送信バッファにすることはできません。なお本節では送信バッファで説明しますが、受信バッファの場合も全く同じです。

一方計算領域を1次元方向に分割した図3-5-4(3)では、着色した部分がメモリー上で不連続なので、着色部分のデータを一度連続な送信バッファに詰め込んでから送信する必要があります。ただし、図3-5-4(4)のように計算領域が縮小(4-5-7節参照)されている場合は連続なので、直接送信することができます。

図3-5-4(5)(6)のように境界部分の1行(または1列)のみを通信する場合も、図3-5-4(6)では着色した部分がメモリー上で不連続なので、着色部分のデータを一度連続な送信バッファに詰め込んでから送信する必要があります。この場合は図3-5-4(7)に示すように、計算領域が縮小されていても不連続になります。

→ 2次元目

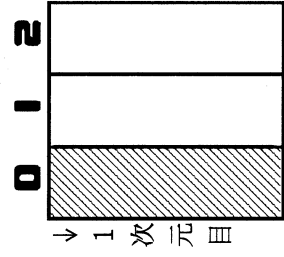


図3-5-4(1) 連続

0 1 2

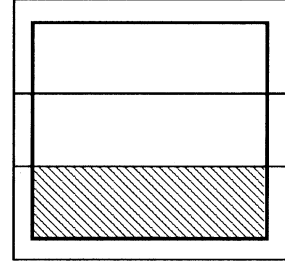


図3-5-4(2) 不連続

0 1 2

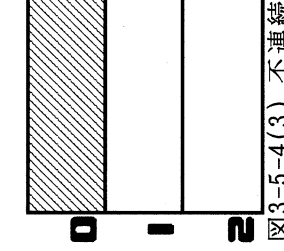


図3-5-4(3) 不連続

0 1 2

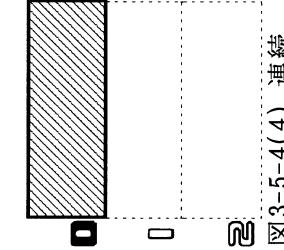


図3-5-4(4) 連続

全体配列

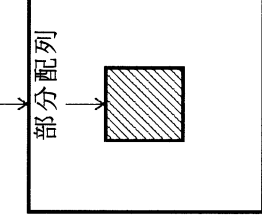


図3-5-5(1)

→ 2次元目

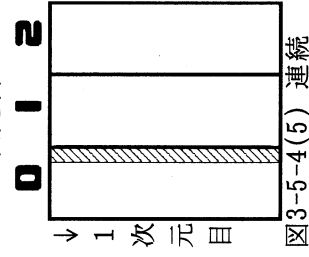


図3-5-4(5) 連続

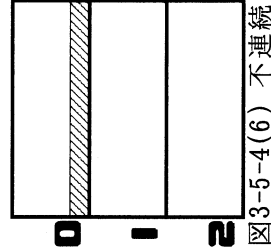


図3-5-4(6) 不連続

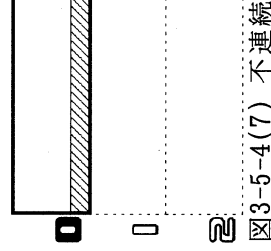


図3-5-4(7) 不連続

全体配列  
部分配列

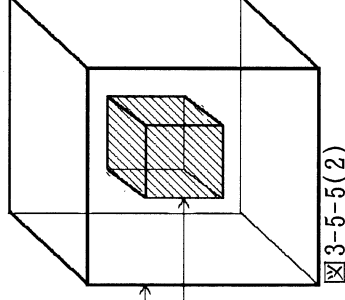


図3-5-5(2)

図3-5-4(2)(3)(6)(7)の着色した部分のようにメモリー上で不連続な部分を送信する場合、送信バッファを新たに確保し、送信バッファに詰め込むロジックをプログラムに追加しなければなりません。また(通常は)受信側でも同様の処理をしなければならず、修正が若干面倒になります。このような場合、派生データ型を使用すると、メモリー上で不連続な部分を直接送受信することができます。

図3-5-4(2)(3)(6)(7)の形状を一般化すると、計算領域が2次元であれば図3-5-5(1)、3次元であれば図3-5-5(2)のようになります。本書では図3-5-5(1)(2)の外枠の配列を全体配列、着色した部分を部分配列と呼びます。部分配列の派生データ型を、前節で紹介したMPIの派生データを作成するサブルーチンで1から作成してもよいのですが、その手間を減らすため、MPI-2ではこのパターンの派生データ型を簡単に作成するサブルーチンMPI\_TYPE\_CREATE\_SUBARRAYが用意されており、使い方を次節で説明します。

一方MPI-2が使用できないマシン環境の場合は、MPI\_TYPE\_CREATE\_SUBARRAYの代わりに、3-5-3-2節で紹介する(筆者が作成した)ユーティリティ・サブルーチンを使用すれば、同様の派生データ型を作成することができます(Fortranのみ)。ただし、MPI\_TYPE\_CREATE\_SUBARRAYは全体配列が何次元でも使用可能なのに対し、ユーティリティ・サブルーチンは2,3次元のみにしか使用できません。

なお、本書の4,5章に、このパターンの派生データ型を使用するプログラム例を数ヶ所掲載していますが、本書のプログラムは基本的にMPI-1の範囲で作成したため、ユーティリティ・サブルーチンを使っています。しかし、MPI-2が使用できる環境の場合は、MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAYを使用するようにして下さい。

MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAYを使用して、図3-5-6(1)に示す全体配列N(2:5,-1:4)内にある、●に示す部分配列の派生データ型を作成する方法を図3-5-6(2)で説明します。なお、MPI\_TYPE\_CREATE\_SUBARRAYの各引数の詳細は付録を参照して下さい。また、MPI-2が使用できないマシン環境の場合は、次節の方法で作成して下さい。

- 図3-5-6(2)の①で、MPI\_TYPE\_CREATE\_SUBARRAYの引数で使用する各配列(整数で名前は任意)を宣言します。図3-5-6(1)の全体配列が2次元なので、配列の大きさは2となります。
- ②で、全体配列の1次元目と2次元目の要素数を配列ISIZEに設定します。
- ③で、部分配列の1次元目と2次元目の要素数を配列ISUBSIZEに設定します。
- ④で、全体配列の各次元の先頭を[i]としたときの、部分配列の先頭の相対座標(図3-5-6(1)の○)を配列ISTARTに指定します。
- ⑤の引数に、全体配列の次元数2、上記で作成した各配列などを指定します。なおC言語の場合、二重線の部分は「MPI\_ORDER\_C」となります。⑤を実行すると、全体配列Nの先頭(★)を起点とした部分配列を表す派生データ型INEWTYPEが作成されます。
- ⑥でINEWTYPEをMPIに登録します。これでINEWTYPEを通信ルーチンの引数で使用することができます。
- ⑦の通信でINEWTYPEを使用します。★の部分を示す配列N(またはN(2,-1))を送受信バッファアの先頭アドレスに指定します。⑦の最初の3つの引数は、送受信バッファアが、「配列Nの先頭から始まる、1個のINEWTYPE型の部分」であることを意味します。

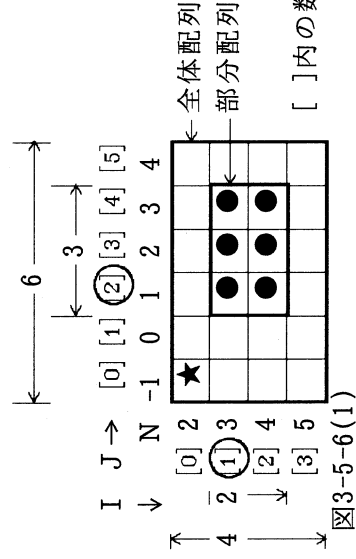


図3-5-6(1) [ ]内の数字は、全体配列の先頭を[i]としたときの相対的な座標

```

:
INCLUDE 'mpif.h'
INTEGER N(2:5, -1:4)
INTEGER ISIZE(2), ISUBSIZE(2), ISTART(2)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
ISIZE(1) = 4
ISIZE(2) = 6
ISUBSIZE(1) = 2
ISUBSIZE(2) = 3
ISTART(1) = 1
ISTART(2) = 2
CALL MPI_TYPE_CREATE_SUBARRAY(2, ISIZE, ISUBSIZE, ISTART,
& MPI_ORDER_FORTRAN, MPI_INTEGER, INEWTYPE, IERR)
CALL MPI_TYPE_COMMIT(INEWTYPE, IERR)
:
CALL MPI_BCAST(N, 1, INEWTYPE, 0, MPI_COMM_WORLD, IERR)
:

```

図3-5-6(2)

### 3-5-3-2 ユーティリティ・サブルーチン

前節で説明したMPI\_TYPE\_CREATE\_SUBARRAYはMPI-2で提供されているサブルーチンなので、MPI-2が使えるいマシン環境では使うことができます。この場合、MPI-1の派生データ型サブルーチンを組み合わせて(筆者が)作成したユーティリティ・サブルーチン(詳細は次ページ以降参照)を使用すれば、部分配列の派生データ型を作成することができます(Fortranのみ)。ただし全体配列が2,3次元の場合にのみ使用可能です。2次元用のサブルーチンPARA\_TYPE\_BLOCK2AとPARA\_TYPE\_BLOCK2の使用例を以下に示します。

#### ■ ユーティリティ・サブルーチンPARA\_TYPE\_BLOCK2A

図3-5-7(1)に示す2次元全体配列N内にある●の部分配列を示す派生データINETYPEを、ユーティリティ・サブルーチンPARA\_TYPE\_BLOCK2A(引数の詳細は次ページ以降参照)を使用して作成するプログラム例を図3-5-7(2)に示します。①の各引数で図3-5-7(1)の○の部分指定します。

PARA\_TYPE\_BLOCK2Aでは、図3-5-7(1)に示すように、作成した派生データ型INETYPEの起点は部分配列の先頭★なので(前節のMPI\_TYPE\_CREATE\_SUBARRAYと違う)、派生データINETYPEを使って通信を行うときは、②の二重線に示すように、★の部分を示すN(3,1)を送受信バッファの先頭アドレスに指定します。②の最初の3つの引数は、送受信バッファが、「配列N(3,1)から始まる、1個のINETYPE型の部分」であることを意味します。

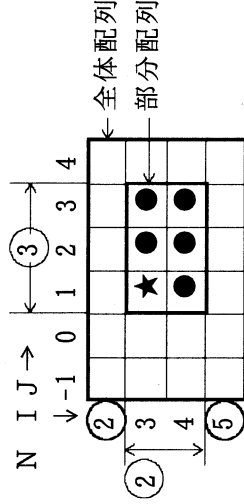


図3-5-7(1)

```

:
INTEGER N(2:5, -1:4)
:
CALL PARA_TYPE_BLOCK2A(2,5,2,3,MPI_INTEGER,INETYPE) ①
CALL MPI_BCAST(N(3,1),1,INETYPE,0,MPI_COMM_WORLD,IERR) ②
:

```

図3-5-7(2)

#### ■ ユーティリティ・サブルーチンPARA\_TYPE\_BLOCK2

一方ユーティリティ・サブルーチンMPI\_TYPE\_BLOCK2では、図3-5-8(2)の③の各引数で図3-5-8(1)の○の部分指定します。図3-5-8(1)に示すように、作成した派生データ型INETYPEの起点は全体配列Nの先頭★なので(前節のMPI\_TYPE\_CREATE\_SUBARRAYと同じ)、派生データINETYPEを使って通信を行うときは、④の二重線に示すように、★の部分を示す配列N(またはN(2,-1))を送受信バッファの先頭アドレスに指定します。②の最初の3つの引数は、送受信バッファが、「配列Nの先頭から始まる、1個のINETYPE型の部分」であることを意味します。

この派生データの先頭から部分配列の先頭までの変位の情報も含むため、作成するときは①より引数が多くなりますが、通信時の送受信バッファの指定は②より簡単になります。

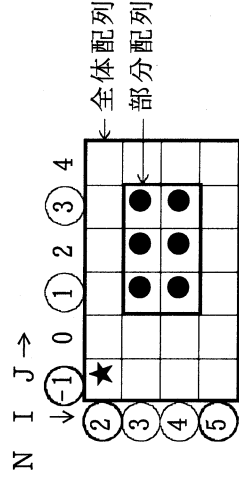


図3-5-8(1)

```

:
INTEGER N(2:5, -1:4)
:
CALL PARA_TYPE_BLOCK2(2,5,-1,3,4,1,3,MPI_INTEGER,INETYPE) ③
CALL MPI_BCAST(N,1,INETYPE,0,MPI_COMM_WORLD,IERR) ④
:

```

図3-5-8(2)

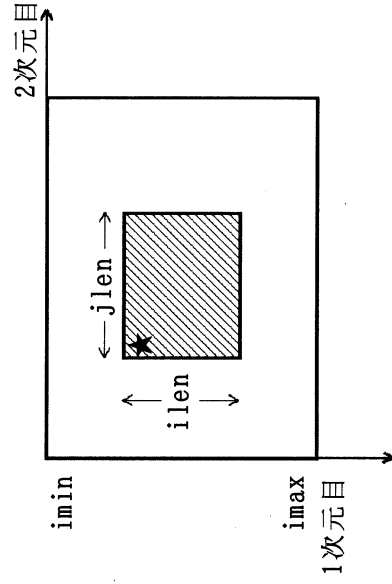
上記の各サブルーチンの3次元版がPARA\_TYPE\_BLOCK3AとPARA\_TYPE\_BLOCK3です。各ユーティリティ・サブルーチンの引数の詳細と実際のプログラムの詳細は次ページ以降を参照して下さい。



(注) 図3-5-9(2)~(4)の各プログラム内のサブルーチンMPI\_TYPE\_SIZEは、本書の2007年1月1日より前の版ではMPI\_TYPE\_EXTENTになっていましたが、MPI\_TYPE\_SIZEの方が適当なので修正しました(MPI\_TYPE\_EXTENTでも正しい結果が得られます)。

■ サブルーチンPARAMETER\_BLOCK2A

要素のデータ型がioldtypeの2次元全体配列(imin:imax, 任意)中にある、大きさがilen×jlenの2次元部分配列を表す派生データ型をinewtypeに戻します。この派生データ型の起点は★なので、MPIで通信を行う際、★の部分を送受信バッファの先頭アドレスに指定して下さい。



CALL PARAMETER\_BLOCK2A(imin, imax, ilen, jlen, ioldtype, inewtype)

- imin : 整数。2次元全体配列の1次元目の下限値を指定します。
- imax : 整数。2次元全体配列の1次元目の上限値を指定します。
- ilen : 整数。2次元部分配列の1次元目の大きさを指定します。
- jlen : 整数。2次元部分配列の2次元目の大きさを指定します。
- ioldtype : 整数。古いデータ型を指定します。例えば整数配列ならばMPI\_INTEGERを指定して下さい。
- inewtype : 整数。作成された新しい派生データ型が戻ります。

```

SUBROUTINE PARAMETER_BLOCK2A
&      (IMIN, IMAX, ILEN, JLEN, IOLDTYPE, INEWTYPE)
  INCLUDE 'mpif.h'
  CALL MPI_TYPE_VECTOR(JLEN, ILEN, IMAX-IMIN+1, IOLDTYPE, INEWTYPE, IERR)
  CALL MPI_TYPE_COMMIT(INEWTYPE, IERR)
END
    
```

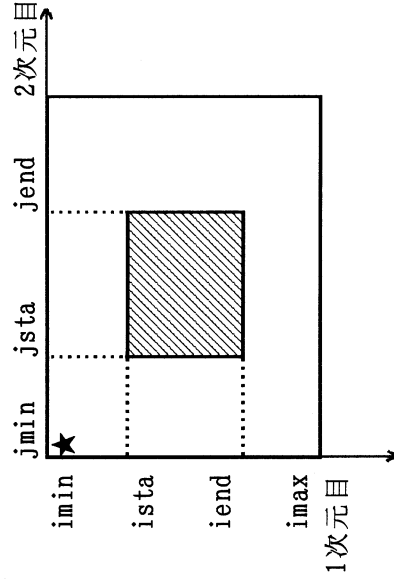
図3-5-9(1)

## ■ サブルーチン PARA\_TYPE\_BLOCK2

(注) 図3-5-9(2)内のMPI\_TYPE\_SIZEは、本書の2007年1月1日より前の版ではMPI\_TYPE\_EXTENTになっていたが、MPI\_TYPE\_SIZEの方が適当なので修正しました(MPI\_TYPE\_EXTENTでも正しい結果が得られます)。

要素のデータ型がioldtypeの2次元全体配列(imin:imax, jmin:jmax)中にある、位置と大きさが

(ista:iend, jsta:jend)の部分配列を表す派生データ型をnewtypeに戻します。この派生データ型の起点は★なので、MPIで通信を行う際、★の部分を送受信バッファの先頭アドレスに指定して下さい。



CALL PARA\_TYPE\_BLOCK2(imin, imax, jmin, ista, iend, jsta, jend, ioldtype, newtype)

- imin : 整数。2次元全体配列の1次元目の下限値を指定します。
- imax : 整数。2次元全体配列の1次元目の上限値を指定します。
- jmin : 整数。2次元全体配列の2次元目の下限値を指定します。
- ista : 整数。2次元部分配列の1次元目の下限値を指定します。
- iend : 整数。2次元部分配列の1次元目の上限値を指定します。
- jsta : 整数。2次元部分配列の2次元目の下限値を指定します。
- jend : 整数。2次元部分配列の2次元目の上限値を指定します。
- ioldtype : 整数。古いデータ型を指定します。例えば整数配列ならばMPI\_INTEGERを指定して下さい。
- newtype : 整数。作成された新しい派生データ型が戻ります。

SUBROUTINE PARA\_TYPE\_BLOCK2

```
& (IMIN, IMAX, JMIN, IJSTA, IJEND, JSTA, JEND, IOLDTYPE, INEWTYPE)
```

```
INCLUDE 'mpif.h'
```

```
INTEGER IBLOCK(2), IDISP(2), ITYPE(2)
```

```
CALL MPI_TYPE_SIZE(IOLDTYPE, ISIZE, IERR)
```

```
ILEN = IEND - IJSTA + 1
```

```
JLEN = JEND - JSTA + 1
```

```
CALL MPI_TYPE_VECTOR(JLEN, ILEN, IMAX-IJMIN+1, IOLDTYPE, ITEMP, IERR)
```

```
IBLOCK(1) = 1
```

```
IBLOCK(2) = 1
```

```
IDISP(1) = 0
```

```
IDISP(2) = ((IMAX-IJMIN+1)*(JSTA-JMIN)+(IJSTA-IJMIN))*ISIZE
```

```
ITYPE(1) = MPI_LB
```

```
ITYPE(2) = ITEMP
```

```
CALL MPI_TYPE_STRUCT(2, IBLOCK, IDISP, ITEMP, INEWTYPE, IERR)
```

```
CALL MPI_TYPE_COMMIT(INEWTYPE, IERR)
```

```
END
```

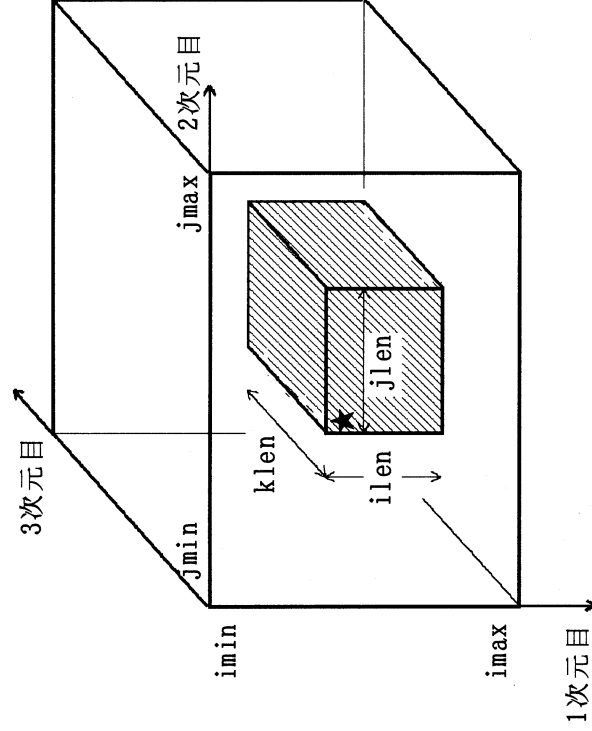
図3-5-9(2)

## ■ サブルーチン PARA\_TYPE\_BLOCK3A

(注) 図3-5-9(3)内のMPI\_TYPE\_SIZEは、本書の2007年1月1日より前の版ではMPI\_TYPE\_EXTENTになっていましたが、MPI\_TYPE\_SIZEの方が適当なので修正しました(MPI\_TYPE\_EXTENTでも正しい結果が得られます)。

要素のデータ型がioldtypeの3次元全体配列(iimin:imax, jmin:jmax, 任意)中にある、大きさが

ilen×jlen×klen の3次元部分配列を表す派生データ型をinewtypeに戻します。この派生データ型の起点は★なので、MPIで通信を行う際、★の部分を送受信バッファアの先頭アドレスに指定して下さい。



```
CALL PARA_TYPE_BLOCK3A(imin, imax, jmin, jmax, ilen, jlen, klen, ioldtype, inewtype)
```

- imin : 整数。3次元全体配列の1次元目の下限値を指定します。
- imax : 整数。3次元全体配列の1次元目の上限値を指定します。
- jmin : 整数。3次元全体配列の2次元目の下限値を指定します。
- jmax : 整数。3次元全体配列の2次元目の上限値を指定します。
- ilen : 整数。3次元部分配列の1次元目の大きさを指定します。
- jlen : 整数。3次元部分配列の2次元目の大きさを指定します。
- klen : 整数。3次元部分配列の3次元目の大きさを指定します。
- ioldtype : 整数。古いデータ型を指定します。例えば整数配列ならばMPI\_INTEGERを指定して下さい。
- inewtype : 整数。作成された新しい派生データ型が戻ります。

```
SUBROUTINE PARA_TYPE_BLOCK3A
```

```
& (IMIN, IMAX, JMIN, JMAX, ILEN, JLEN, KLEN, IOLDTYPE, INEWTYPE)
```

```
INCLUDE 'mpif.h'
```

```
CALL MPI_TYPE_SIZE(IOLDTYPE, ISIZE, IERR)
```

```
CALL MPI_TYPE_VECTOR(JLEN, ILEN, IMAX-IMIN+1, IOLDTYPE, ITEMP, IERR)
```

```
IDIST = (IMAX-IMIN+1)*(JMAX-JMIN+1)*ISIZE
```

```
CALL MPI_TYPE_HVECTOR(KLEN, 1, IDIST, ITEMP, INEWTYPE, IERR)
```

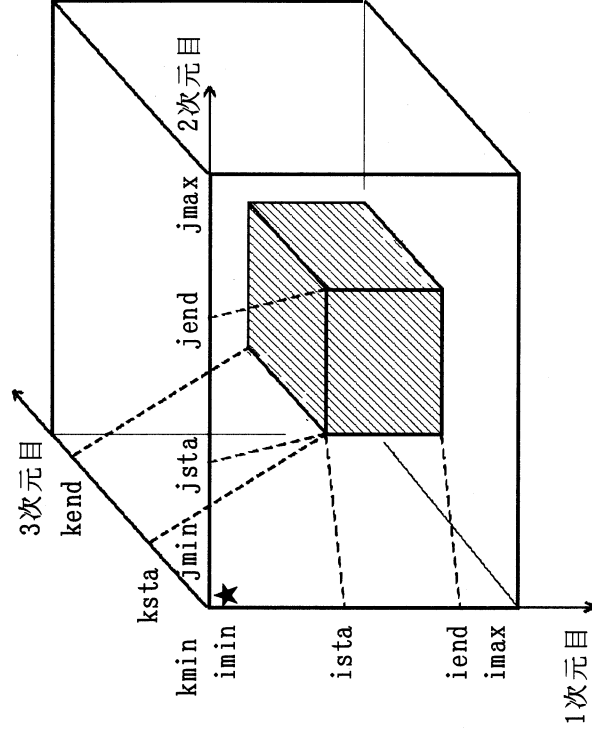
```
CALL MPI_TYPE_COMMIT(INEWTYPE, IERR)
```

```
END
```

図3-5-9(3)

(注) 図3-5-9(4)内のMPI\_TYPE\_SIZEは、本書の2007年1月1日より前の版ではMPI\_TYPE\_EXTENTになっていましたが、MPI\_TYPE\_SIZEの方が適当なので修正しました(MPI\_TYPE\_EXTENTでも正しい結果が得られます)。

要素のデータ型がioldtypeの3次元全体配列(imin:imax, jmin:jmax, kmin:kmax)中にある、位置と大きさが(ista:iend, jsta:jend, ksta:kend)の3次元部分配列を表す派生データ型をinewtypeに戻します。この派生データ型の起点は★なので、MPIで通信を行う際、★の部分を送受信バッファの先頭アドレスに指定して下さい。



CALL PARA\_TYPE\_BLOCK3(imin, imax, jmin, jmax, kmin, kmax, ista, iend, jsta, jend, ksta, kend, ioldtype, inewtype)

- imin : 整数。3次元全体配列の1次元目の下限値を指定します。
- imax : 整数。3次元全体配列の1次元目の上限値を指定します。
- jmin : 整数。3次元全体配列の2次元目の下限値を指定します。
- jmax : 整数。3次元全体配列の2次元目の上限値を指定します。
- kmin : 整数。3次元全体配列の3次元目の下限値を指定します。
- kmax : 整数。3次元全体配列の3次元目の上限値を指定します。
- ista : 整数。3次元部分配列の1次元目の下限値を指定します。
- iend : 整数。3次元部分配列の1次元目の上限値を指定します。
- jsta : 整数。3次元部分配列の2次元目の下限値を指定します。
- jend : 整数。3次元部分配列の2次元目の上限値を指定します。
- ksta : 整数。3次元部分配列の3次元目の下限値を指定します。
- kend : 整数。3次元部分配列の3次元目の上限値を指定します。
- ioldtype : 整数。古いデータ型を指定します。例えば整数配列ならばMPI\_INTEGERを指定して下さい。
- inewtype : 整数。作成された新しい派生データ型が戻ります。

```

SUBROUTINE PARA_TYPE_BLOCK3(IMIN, IMAX, JMIN, JMAX, KMIN,
&   ISTA, IEND, JSTA, JEND, KSTA, KEND, IOLDTYPE, INEWTYPE)
  INCLUDE 'mpif.h'
  INTEGER IBLOCK(2), IDISP(2), ITYPE(2)
  CALL MPI_TYPE_SIZE(IOLDTYPE, ISIZE, IERR)
  ILEN = IEND - ISTA + 1
  JLEN = JEND - JSTA + 1
  KLEN = KEND - KSTA + 1
  CALL MPI_TYPE_VECTOR(JLEN, ILEN, IMAX-IMIN+1, IOLDTYPE, ITEMP, IERR)
  IDIST = (IMAX-IMIN+1)*(JMAX-JMIN+1)*ISIZE
  CALL MPI_TYPE_HVECTOR(KLEN, 1, IDIST, ITEMP, ITEMP2, IERR)
  IBLOCK(1) = 1
  IBLOCK(2) = 1
  IDISP(1) = 0
  IDISP(2) = ( (IMAX-IMIN+1)*(JMAX-JMIN+1)*(KSTA-KMIN)
&   + (IMAX-IMIN+1)*(JSTA-JMIN) + (ISTA-IMIN) ) * ISIZE
  ITYPE(1) = MPI_LB
  ITYPE(2) = ITEMP2
  CALL MPI_TYPE_STRUCT(2, IBLOCK, IDISP, ITYPE, INEWTYPE, IERR)
  CALL MPI_TYPE_COMMIT(INEWTYPE, IERR)
END

```

図3-5-9(4)

### 3-5-4 構造体を表す派生データ型の作成

構造体で定義した変数や配列は、MPIが提供するデータ型(MPI\_INTEGERなど)では通信できないので、構造体を表す派生データ型を作成して通信を行います。図3-5-10(1)の構造体の派生データ型を作成するプログラム例を図3-5-10(2)に示します。

- ①で図3-5-10(1)に示すKOUZOUTAIという名の構造体を定義し、②で変数WORKをこの構造体で宣言します。
- ③～⑥で、この変数WORKを通信するために使用する派生データ型IKOUZOUを、MPIのサブルーチンMPI\_TYPE\_STRUCTを使用して作成します。MPI\_TYPE\_STRUCTの引数の詳細は付録を参照して下さい。図3-5-10(1)の構造体KOUZOUTAIは、以下の3つの部分(以後ブロックと呼びます)から構成されています。まず③で、MPI\_TYPE\_STRUCTの各引数に使用する、大きさ3(=ブロック数)の配列(整数で名前は任意)を宣言します。

```
[ブロック1] 4バイトの整数が2個
[ブロック2] 8バイトの倍精度実数が1個
[ブロック3] 4バイトの単精度実数が2個
```

- ④で配列IBLOCKに、上記の各ブロックに含まれる変数の数(上記下線部の数字)を代入します。
- ⑤で配列IDISPに、構造体の先頭から各ブロックの先頭までの変位(バイト)を代入します(図3-5-10(1)の○参照)。
- ⑥で配列IATYPEに、上記の各ブロックに含まれる変数のデータ型を指定します。
- ⑦でMPIのサブルーチンMPI\_TYPE\_STRUCTをコールすると、下線部のIKOUZOU(任意の整数)に、構造体KOUZOUTAIを表す派生データ型が戻ります。

- ⑧で派生データ型IKOUZOUをMPIに登録します。これでIKOUZOUをMPIの通信のデータ型として使うことができます。

- ⑨でランク0は、構造体で宣言した変数WORKに値を代入し、⑩で派生データ型IKOUZOUを使用して他のプロセスに通信します。言うまでもありませんが、⑩で例えば配列WORK(2)の2つのデータを送信する場合は「CALL MPI\_BCAST(WORK,2,IKOUZOU,~)」のようになります。

```
先頭からの [ブロック1] [ブロック2] [ブロック3]
変位(バイト) ① 4 ⑧ ⑩ ⑩
KOUZOUTAI I J X X A B
```

整数 倍精度実数 単精度実数

図3-5-10(1)

```
TYPE:: KOUZOUTAI
INTEGER:: I
INTEGER:: J
REAL*8 :: X
REAL :: A
REAL :: B
END TYPE KOUZOUTAI
TYPE(KOUZOUTAI) WORK
INTEGER IBLOCK(3), IDISP(3), IATYPE(3)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
```

図3-5-10(2)

```
IBLOCK(1) = 2
IBLOCK(2) = 1
IBLOCK(3) = 2
IDISP (1) = 0
IDISP (2) = 8
IDISP (3) = 16
IATYPE (1) = MPI_INTEGER
IATYPE (2) = MPI_REAL8
IATYPE (3) = MPI_REAL
CALL MPI_TYPE_STRUCT(3,IBLOCK, IDISP, IATYPE,
& IKOUZOU, IERR)
CALL MPI_TYPE_COMMIT(IKOUZOU, IERR)
IF (MYRANK==0) THEN
  WORK%I = 1
  WORK%J = 2
  WORK%X = 3.0D0
  WORK%A = 4.0
  WORK%B = 5.0
ENDIF
CALL MPI_BCAST(WORK,1,IKOUZOU, 0,
& MPI_COMM_WORLD, IERR)
:
```

3-2節で説明したように、プロセスの集合からなるグループのグループ名をコミュニケーションと呼びます。このうち、全プロセスから構成されるグループはMPIがあらかじめ作成しており、そのコミュニケーションをMPI\_COMM\_WORLDと呼びます。これとは別に、任意のプロセスから構成される新グループをユーザーが作成することができま

3-1節で紹介したMPIサブルーチンの分類のうち、コミュニケーションに所属するサブルーチンで新グループを作成しますが、このうちMPI\_COMM\_SPLITのみでも新グループを作成できるので、本書ではMPI\_COMM\_SPLITのみを紹介します。

付録のMPI\_COMM\_SPLITの説明を一読して下さい。...

某国の政界では、新グループや新党、新党が頻繁に発生と消滅を繰り返していますが(本書の初版を作成した1995年頃の話です)、実はMPIの世界では、新グループを作成する必要がある場合はほとんどありません。しいて挙げれば、全く別の複数のプログラム(例えば構造解析と流体解析)をそれぞれ並列化して実行する、連成解析のような場合に必要となります(4-7-1-2節参照)。

上記で紹介したMPI\_COMM\_SPLITを使用して、例えば図3-6-1に示すように、構造解析を担当するプロセスのグループ(コミュニケーションIKOUZOU)と、流体解析を担当するプロセスのグループ(コミュニケーションIRYUTAI)を作成します。このときグループ内の各プロセスには、ランクの値を付けます(下記の例ではランク0, 1, 2)。

例えば構造解析のプロセス間で通信を行う場合、通信ルーチンの引数にはコミュニケーションIKOUZOUを指定します。一方、構造解析のプロセスと流体解析のプロセス間で通信を行う場合、通信ルーチンの引数にはMPI\_COMM\_WORLDを指定します。また通信ルーチンの引数にランクの値を指定する場合、コミュニケーションで指定したグループ内でのランクの値を指定します。

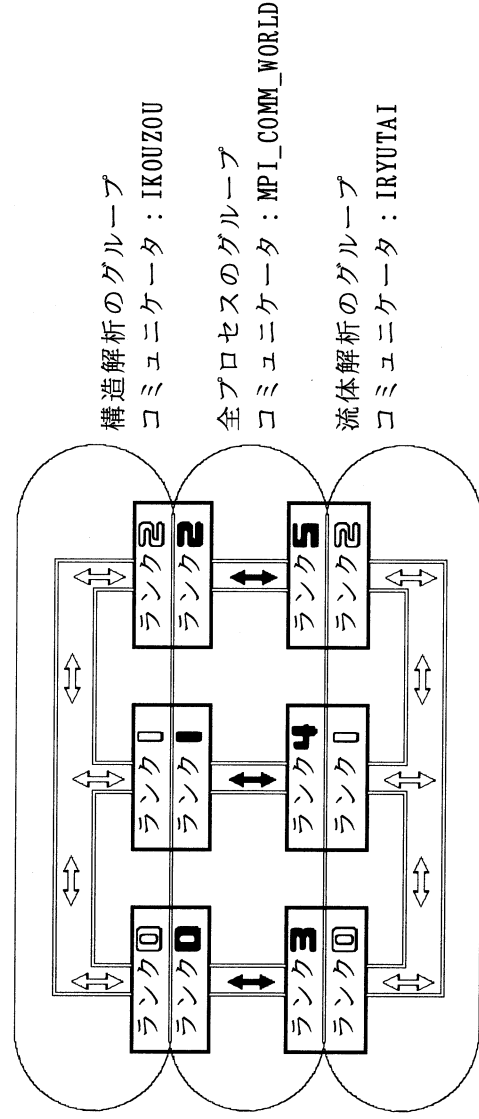


図3-6-1

## 3-7 C言語のMPI関数

C言語でもMPIの関数を使用することができます。以下に、本書の付録に掲載したMPI\_BCASTのサンプルプログラムをC言語で書いたプログラムを示します。

```
#include "mpi.h"
int main(argc, argv)
int argc;
char **argv;
{
    int nprocs, myrank, ibuf;
    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &nprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
    if (myrank == 0) {
        ibuf = 1;
    } else {
        ibuf = 0;
    }
    MPI_Bcast(&ibuf, 1, MPI_INT, 0, MPI_COMM_WORLD);
    printf("MYRANK = %d IBUF = %d\n", myrank, ibuf);
    MPI_Finalize();
}
```

図3-7-1

FortranとC言語のMPIの主な相違点を以下に示します。

なお、本書の付録に掲載したのはFortran版のMPIルーチンの使用方法です。例えばFortran版で引数が整数になっていても、C言語版では別のデータ型で宣言しなければなりませんので、C言語のユーザの方は、参考文献[15][27]のC言語の説明箇所を参照して下さい。またC言語用のMPIのテキストが参考文献[37][38]にあります。

●C言語のMPI関数は大文字と小文字の区別があり、関数名の一部と、MPIが使用する定数(MPI\_COMM\_WORLDなどは大文字になります)。

●引数はFortranとC言語でほぼ同じですが、Fortranの最後の引数*ierror*(戻り値が返る)がC言語にはなく、下記に示すようにMPIの関数自体が戻り値となります。

```
int ierr;
ierr = MPI_Finalize();
```

●引数の一部がポインタになります。

●Fortranのインクルードファイル「mpif.h」は、C言語では「mpi.h」となります。

●例えばFortranの「INTEGER ISTATUS(MPI\_STATUS\_SIZE)」の宣言は、C言語では「MPI\_Status istatus;」となります。

●MPIの通信ルーチンの送受信バッファに多次元配列(のポインタ)を指定する場合、配列内の各要素がメモリー上で連続している多次元配列のみ、指定することができます。



C言語の場合のMPIの基本データ型、MPIの定義済み演算を以下に示します。MPI\_FLOAT\_INTなど(FortranのMPI\_2REALなどに相当)は、異なる2つのデータ型のペアを(構造体を使用して)扱うことができます(3-3-6節参照)。

	C言語のデータ型	MPIの基本データ型
整数型	signed short int signed int signed long int unsigned short int unsigned int unsigned long int (longlong int)	MPI_SHORT MPI_INT MPI_LONG MPI_UNSIGNED_SHORT MPI_UNSIGNED MPI_UNSIGNED_LONG (MPI_LONG_LONG_INT)
実数型	float double long double	MPI_FLOAT MPI_DOUBLE MPI_LONG_DOUBLE
その他	signed char unsigned char 対応なし	MPI_CHAR MPI_UNSIGNED_CHAR MPI_BYTE, MPI_PACKED

図3-7-2(1)

(注)

マシン環境によっては  
カッコをつけたデータ型などが  
提供されている場合があります。

	MPIの基本データ型	MPIの定義済み演算
MPIの整数型と実数型		MPI_SUM (合計) MPI_PROD (積)
上記の整数型と実数型		MPI_MAX (最大) MPI_MIN (最小)
MPI_FLOAT_INT MPI_DOUBLE_INT MPI_LONG_INT MPI_2INT MPI_SHORT_INT MPI_LONG_DOUBLE_INT		MPI_MAXLOC (最大と位置) MPI_MINLOC (最小と位置)
上記の整数型		MPI_LAND (論理AND) MPI_LOR (論理OR) MPI_LXOR (論理XOR)
上記の整数型とMPI_BYTE		MPI_BAND (ビットAND) MPI_BOR (ビットOR) MPI_BXOR (ビットXOR)

図3-7-2(2)

ただし、

```

MPI_FLOAT_INT = {MPI_FLOAT, MPI_INT}
MPI_DOUBLE_INT = {MPI_DOUBLE, MPI_INT}
MPI_LONG_INT = {MPI_LONG, MPI_INT}
MPI_2INT = {MPI_INT, MPI_INT}
MPI_SHORT_INT = {MPI_SHORT, MPI_INT}
MPI_LONG_DOUBLE_INT = {MPI_LONG_DOUBLE, MPI_INT}
    
```

## 3-8 MPI-2

3-1節で紹介したように、MPIの新機能がMPI-2で追加になりました。本節ではそのうちの一部を紹介いたします。これらの機能は、使用するマシン環境にMPI-2が導入されている場合のみ、使用することができます。

### 3-8-1 MPI\_IN\_PLACE

3-3-5節で説明したように、下記の集団通信ルーチンには、「送信バッファと受信バッファ(の実に使用する部分)はメモリー上で重なってはならない」という制限があります。このため通信の際、送信バッファとは別に、送信バッファと同じ大きさの受信バッファを新たに用意しなければならず、バッファの大きさが大きいときはメモリーが無駄になりました。

- (グループ2) **○**MPI\_GATHER, **○**MPI\_SCATTER, **○**MPI\_ALLGATHER, **✕**MPI\_ALLTOALL
- MPI\_GATHERV, **○**MPI\_SCATTERV, **✕**MPI\_ALLGATHERV, **✕**MPI\_ALLTOALLV, **✕**MPI\_ALLTOALLV
- (グループ3) **○**MPI\_REDUCE, **○**MPI\_ALLREDUCE, **○**MPI\_SCAN, **✕**MPI\_EXSCAN, **○**MPI\_REDUCE\_SCATTER

MPI-2では、上記のうち**○**をつけたルーチンについては、以下に示す方法を使用すれば「送信バッファと受信バッファをメモリー上で同一にすることができ

ます。まず図3-8-1(1)に今までの方法を示します。送信バッファSにセットされた送信データが、MPI\_ALLREDUCEによって図3-8-1(2)に示すように通信しながら合計され、結果が全プロセスの受信バッファSSに入ります。上記の制限があるため、送信バッファSとは別に、受信バッファSSを用意しなければなりません。

一方MPI-2の新機能を使用した図3-8-2(1)では、(本来送信バッファを指定する)1つ目の引数に、MPI-2で新たに提供された変数MPI\_IN\_PLACEを指定します。そして(本来受信バッファを指定する)2つ目の引数に指定した変数は、送信バッファを兼用します。まず図3-8-2(2)に示すように変数Sに送信データをセットし、MPI\_ALLREDUCEをコールします。すると通信しながら合計された結果が変数Sにセットされます(元の値が更新されます)。このため受信バッファSSを新たに用意する必要はありません。なお、MPI\_IN\_PLACEの使い方は各集団通信サブルーチンで若干異なりますので、付録の「MPI\_IN\_PLACEの使い方」を参照して下さい。

この機能は、使用するMPIがMPI-2のMPI\_IN\_PLACEオプションをサポートしている場合にのみ使用することができます。MPI-2をサポートしているMPIでも、MPI-2の全ての機能はサポートしておらず、MPI\_IN\_PLACEを使用することができない場合もありますので、MPI\_IN\_PLACEオプションが使用できるかどうかを、マニュアルなどで確認してからお使い下さい。

```

:
CALL MPI_ALLREDUCE
& (S,SS,1,MPI_INTEGER,
& MPI_SUM,MPI_COMM_WORLD,IERR)
S = SS
PRINT *,S
:

```

図3-8-1(1)

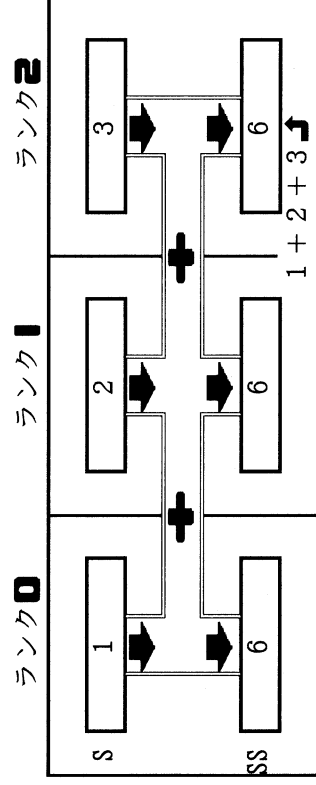


図3-8-1(2)

```

:
CALL MPI_ALLREDUCE
& (MPI_IN_PLACE,S,1,MPI_INTEGER,
& MPI_SUM,MPI_COMM_WORLD,IERR)
PRINT *,S
:

```

図3-8-2(1)

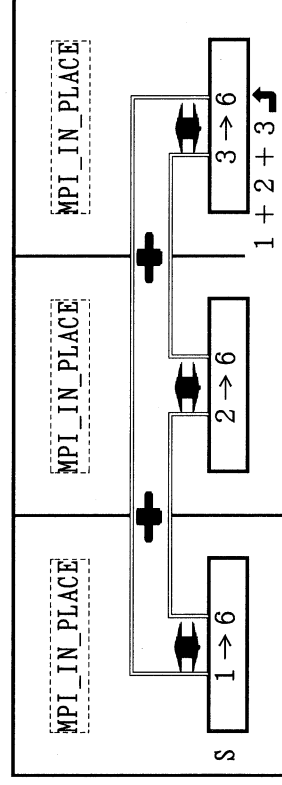


図3-8-2(2)

## 第4章

# プログラムの並列化の方法

第3章で、主なメッセージ交換サブルーチンの使用方法について説明しました。英語で言えばこれらで主な英単語を覚えたことになります。ここからは、覚えた単語を使って会話を行うフェーズになります。本章では、実際のプログラムを並列化する方法について説明します。

## 4-1 パフォーマンスに関する考慮点

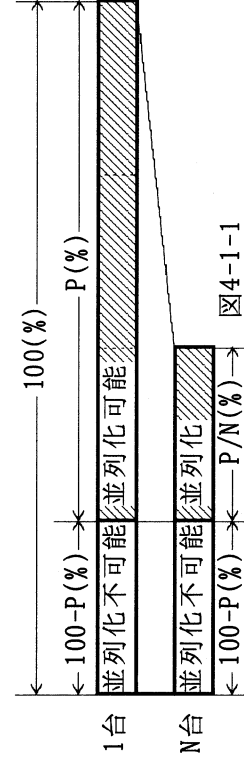
例えばノード数8台で並列プログラムを実行した場合、1台の場合の8倍とは言わないまでも、7倍程度の速度向上が得られるのではないかと考えがちです。勿論8倍近く出るプログラムもありますが、台数を増やせば増やすほど却って遅くなってしまいうプログラムもあります。この差は一体どこから来るのでしょうか。

### 4-1-1 アムダールの法則

図4-1-1で、単体プログラムを実行した場合の経過時間(の割合)を100(%)とします。このプログラムのP%(本例では75%)の部分が並列化可能とします(これを並列化率がP%)です。N台(本例では3台)で並列に実行すると、並列化されていない100-P(%)の部分はそのまま残り、並列化されたP(%)の部分はN台で実行してP/N(%)となるので、合計(100-P)+P/N(%)となります。1台の経過時間と比べてN台で何倍速くなったかを表す比率を速度向上率と呼び、本例では3台で並列に実行して速度向上率は2倍となります。これを一般化すると以下のようになります、アムダールの法則と呼びます。ただし、この式および以下の説明では、並列化に伴う通信などのオーバーヘッドは考慮に入っていません。

$$\text{並列化率P(\%)} \text{のプログラムをノード数N台で実行したとき、}$$

$$\text{速度向上率(倍)} = \frac{100}{(100-P)+P/N} \dots 1 \text{台での経過時間(}\% \text{)}$$



P: 並列化率(%)	50	60	(70)	80	90	(95)	100
ノード数 2台	1.33	1.43	1.54	1.67	1.82	1.90	2.0
ノード数 8台	1.78	2.1	(2.58)	3.33	4.7	(5.93)	8.0
ノード数 ∞台	2.0	2.5	3.33	5.0	10.0	20.0	∞

図4-1-2(1)

ノード数が2台、8台、∞台の場合の速度向上率を図4-1-2(1)に、グラフを図4-1-2(2)に示します。

●ノード数が8台の(多い)場合と2台の(少ない)場合を比較してみます。8台の場合、図4-1-2(1)から分かるように、並列化率が70%程度だと速度向上率は2.58倍にしかならず、95%でようやく5.93倍になります。つまりノード数が多い場合、並列化率がかなり高くないと並列化の効果は得られません。

一方ノード数が2台の場合、並列化率が70%程度でも1.54倍となります。つまりノード数が少ない場合、並列化率が多少低くても、ある程度並列化の効果は得られます。

別の見方として、並列化率が70%の場合、ノードを1台から1台増やしただけで1.54倍になりますが、そこからさらに6台増やしても2.58倍にしかならないと言えます。

●アムダールの法則でノード数Nを∞台にすると、図4-1-2(3)に示すように並列化可能な部分はほぼゼロになります(ただし実際は通信などのオーバーヘッドが加わります)、並列化不可能な部分が残ってしまいます。例えば並列化率が80%の場合、台数をいくら増やしても速度向上率は5倍以上にはなりません。すなわち図4-1-2(1)(2)で∞台の場合の速度向上率は、各並列化率での速度向上率の上限を意味します。

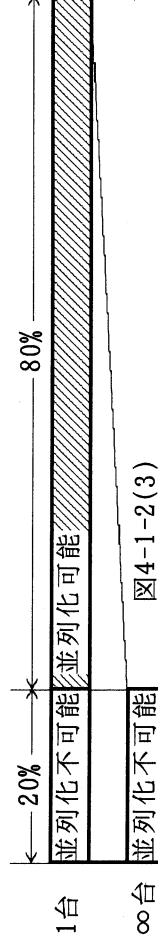


図4-1-2(3)

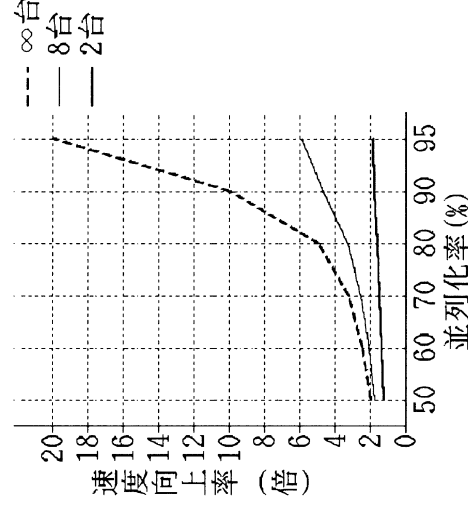


図4-1-2(2)

## 4-1-2 パフォーマンスを向上させるための考慮点

速度向上率についてもう少し細かく検討してみます。例えば図4-1-3(1)のようにノード数が4台の場合を考えます。並列化可能な80%の部分が並列化されますが、この例では $80 \div 4$ の20%にならず、1台目のノードは15%、4台目のノードは25%などと、計算時間がバラついてしまっています。並列ジョブでは、このようにバラつきがある場合、一番遅いプロセスに足をひっぱられてしまいます。言い換えると、処理を行わずに遊んでいるプロセスがあるので、計算時間が遅くなります。

これに加え、並列化にともなう通信時間が増加します。その結果、図4-1-3の例では、4台で2倍の速度向上率しか得ることができませんでした。

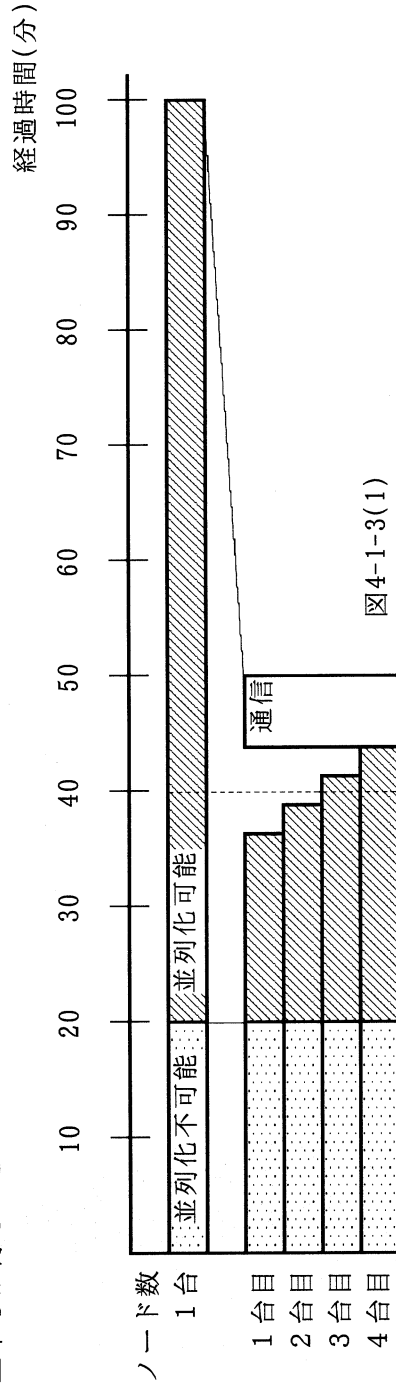


図4-1-3(1)

以上のことから、速度向上率を上げるためには、以下の点に考慮する必要があることが分かります。

- (1) (ノード数が多い場合)並列化率をなるべく高くする。
- (2) プロセス間のロードバランスをなるべく均等ににする。
- (3) 通信時間をなるべく少なくする。

●(1)はプログラムのアルゴリズムの問題であり、本質的に並列性のないプログラムを並列化することは不可能ですが、アルゴリズムによっては並列性を持つアルゴリズムに変更できる場合もあります。

例えば差分法で使用されるSOR法は、そのままでは計算に依存関係があるため並列化できませんが、ゼブラなどの方法に変更すると、依存関係がなくなり並列化できるようになります(5-6節参照)。また、一次逐次演算(4-6-9節)も同様です。

●(2)に関しては、ブロック分割(D0ループの反復を例えばプロセス数が4ならば $1/4$ ずつ分割)で分割すれば、多くの場合ロードバランスは均等になります。ブロック分割でロードバランスが均等にならない(ロードインバランスと言います)例と対処方法については4-5-3節で説明します。

●(3)がもっとも注意すべき項目です。分散メモリー型並列計算機での並列化はまさに通信時間との戦いです。並列化を行うときは通信時間をいかにして短縮するかを常に念頭に置く必要があります。以下、通信時間を短縮するための考慮点について述べます。

●余談ですが、プログラムによっては並列化するとCPU台数以上の速度向上率(例えば3台で3倍以上)が出ることもあり、これをスーパーリニアと呼ぶことがあります。この理由ですが、単体プログラムで図4-1-3(2)の配列内の要素を数字の順に処理するとします。このようにメモリー内の要素を飛び飛びに処理すると、スカラ計算機(ワークステーションやパソコン)の場合はキャッシュミスが発生して速度が低下します(参考文献[6]の4章参照)。これを並列化した図4-1-3(3)でもメモリー内の要素を飛び飛びに処理しますが、1プロセス当たりの担当する範囲が狭くなり、前に処理した要素の近隣の要素がキャッシュに残っている可能性が高くなるため、キャッシュミスによる速度低下の影響が少なくなります。このような現象は、計算領域が非構造格子から構成されていて間接アドレスを使ったプログラムなどで現れます。

### 1 プロセス



図4-1-3(2)

図4-1-3(3)

## ■ 通信の特性

通信の特性は一般に次のようになっています(ただしマシン環境によって異なります)。

- 図4-1-4に示すように、通信時間は通信バイト数に比例します。すなわち、通信バイト数が2倍になれば通信時間も2倍になります。
- 図4-1-4に示すように、一回通信を行うごとに一定の立上り時間(Latency)がかかります。
- 並列計算機内の各ノード間の通信時間はほぼ同じであるマシンが一般的ですが、これはマシン環境に依存します。

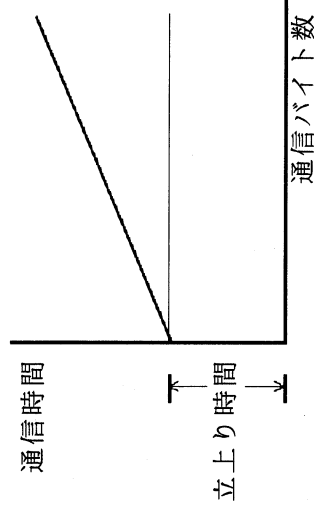


図4-1-4

## ■ 通信に関する考慮点(1)

以上の通信特性から、ユーザーは次の2点に考慮する必要があります。

(1) 同じ通信量なら通信回数を減らす。

例えば要素数が100個のデータを1度に送信すれば、立上り時間は1回で済みますが、これを1要素ずつ100回送信すると、立上り時間も100回分かってしまいますので、なるべく通信回数を少なくして下さい。

では、このような状態はどのような場合に発生するのでしょうか？ 例えば図4-1-5のように、配列Aを1次元目(I方向)で分割しており、ランクRのプロセスは斜線の1行を他のプロセスに送信する必要があるとします。この場合、3-1節「付録について」で説明したように、この配列はメモリー上で図の縦方向に連続しています(図の↓参照)。従って、斜線の部分は図4-1-5では連続して見えますが、実際には各要素はメモリー上でとびとびになっています。

この場合、(a)のように1度に1要素ずつ100回送信すると、100回分の立ち上がり時間がかかってしまいます。これに対して、(b)のように送信したい要素をいったん連続した送信バッファー(配列)に入れ、その後で一気に送信すれば、立ち上がり時間は1回で済むわけです。この場合、もちろん要素を送信バッファーに詰め込むオーバーヘッドはかかりますが、通信というのはI/Oのようなものなので、それに比べるとCPUが行う詰め込むための時間は(一般に)無視することができます。

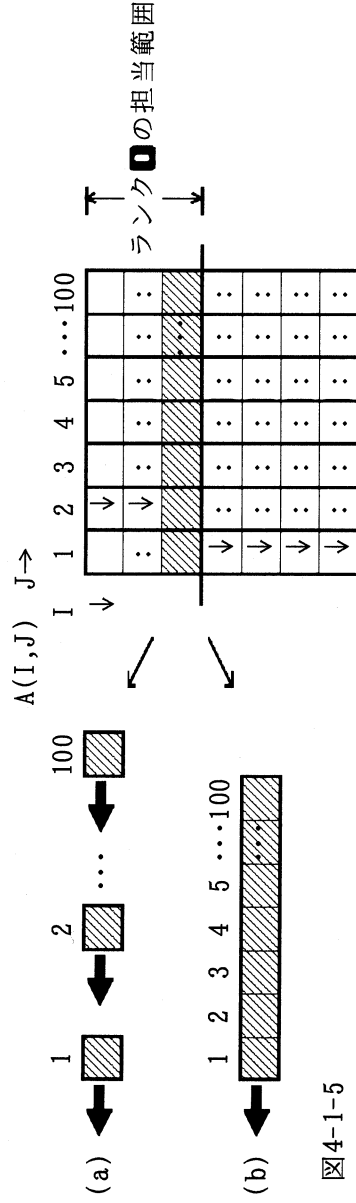


図4-1-5

なお、MPIで提供されている派生データ型(3-5節参照)を使用すれば、図4-1-5のようにメモリー上で規則的に離れて配置されているデータは直接送信できますので、通信回数は1度で、さらに送信バッファ一に詰め込む作業も不要になります(ただしパフォーマンスが良いかどうかはマシン環境に依存します)。

初心者が並列化したプログラムを見ると、図4-1-6(1)のように、D0ループ内でメッセージ交換を行っていることがわかります。この例では10000回「CALL MPI\_BCAST」を実行するため、10000回分の立ち上がり時間がかかってしまいます。従って、メッセージ交換はD0ループ内では行わず、図4-1-6(2)のようにD0ループの外で一度だけ行うのが普通です。

```

:
DO I=1,10000
  計算
  CALL MPI_BCAST(A(I),~)
ENDDO
:

```

図4-1-6(1) ✘

```

:
DO I=1,10000
  計算
ENDDO
CALL MPI_BCAST(A,~)
:

```

図4-1-6(2) ○

■ 通信に関する考慮点(2)

次に2つ目の考慮点を以下に示します。

(2) 通信量自体が少なくなるようにする(配列の分割方法に注意)。

通信を行う必要が生じた場合でも、並列化の方法によっては通信量を少なくできる場合があります。差分法では、あるメッシュの値を、前後左右のメッシュの値から求める演算が発生します。図4-1-7(1)では長方形の計算領域を1次元目で分割しています。ランク0のプロセスが①の計算をする場合、前後左右の●はすべて自分が持っているので問題ありません。ところが②の計算をする場合、■は持っていないが、□はランク1のプロセスが担当するので持っています。従って、計算を行う前に、斜線の1行のデータをランク1のプロセスから送ってもらう必要があります。

この場合、図4-1-7(2)のように計算領域を分割することもできます。両者を比較すると、図4-1-7(2)の方が境界の長さが短いため、通信量も少なくなります。このように、計算領域が長方形(2次元)や直方体(3次元)の場合、どの次元で分割すれば通信量が少なくなるかを考慮に入れる必要があります。これについては4-5-8節で説明します。

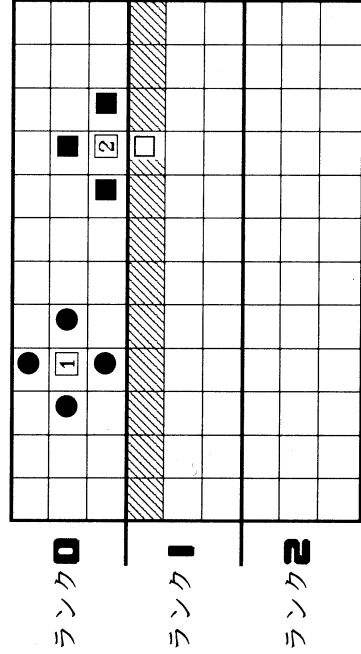


図4-1-7(1)

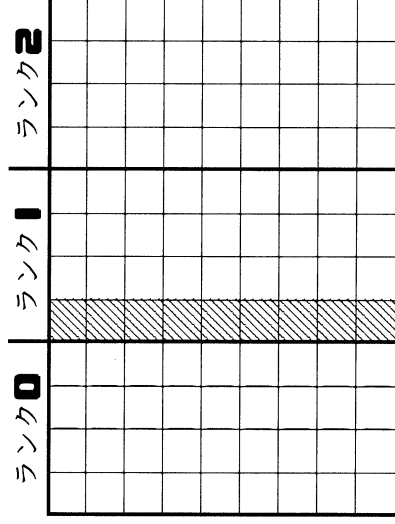


図4-1-7(2)

## ■ 計算量と通信量の関係

並列化の効果を見積もる一つの目安を紹介いたします。図4-1-8(1)で各プロセスは一辺がnのオーダーの領域を計算し、その後自分が計算した部分を他のプロセスに通信いたします。この場合、大ざっぱに言うると計算量は $\alpha n^2$ 、通信量は $\beta n^2$ とどちらも $n^2$ のオーダーになります。計算に比べると、通信というのは1/0のようなものですから、係数は $\alpha \ll \beta$ になります。従って、(通信量)/(計算量) =  $\beta/\alpha \gg 1$ となりますので、図4-1-8(1)のように計算量と通信量が同じ次数の場合、よほど計算部分が多くなくと通信量の割合が多くなって並列化の効果が出ません。

一方、図4-1-8(2)のように、境界の1列のみを通信する場合、通信量は $\beta n$ とオーダーがnになるので、(通信量)/(計算量) =  $\beta/\alpha n$  となり、特にnが大きくなるほど通信量の割合は減少します。

以上のことから、一般に、計算量より通信量が少なくなるとも次数が1つ以上少なくなると通信量を並列化すれば、特に問題のサイズが大きければ通信量の割合が減少し、並列化の効果が出ると言うことができます。

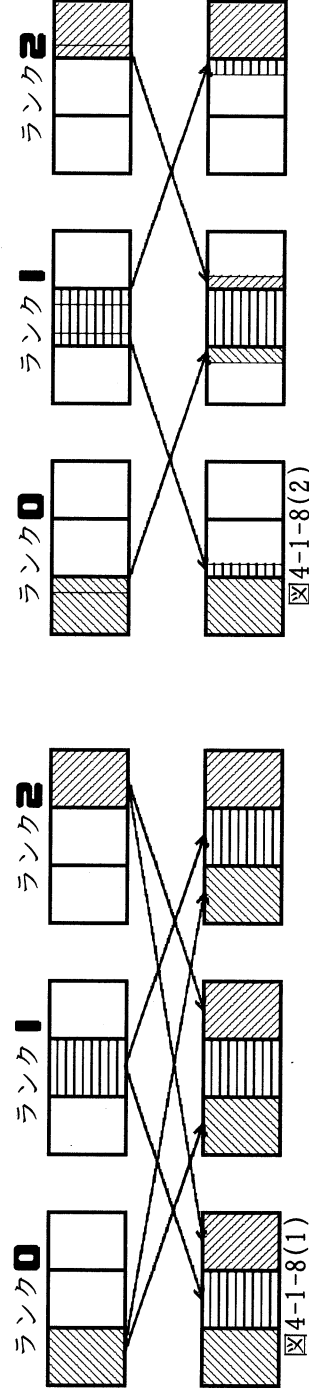


図4-1-8(1)

図4-1-8(2)

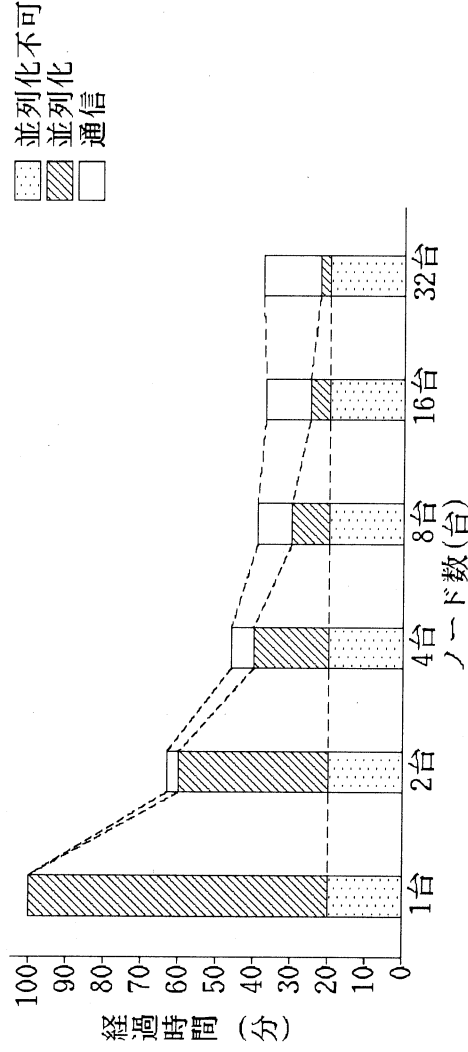
## ■ ノード数を増やすとなぜサチュレートするのか

大部分の並列プログラムでは、ノード数を増やしていくとどこかでサチュレートし、それ以上はノードを増やすとかえって遅くなってしまいます。なぜこのようになるのでしょうか？ これを以下のグラフで示します。なお、この例は図4-1-2と同じく並列化可能部分が80%の例です。

まず並列化できない部分に関しては、ノード数を何台にしようとも、20%のまま残ります。次に並列化可能部分ですが、これはプロセス間のロードバランスが均一であれば、2台でほぼ1/2、4台なら1/4、... のように減少していきます。

問題なのは通信時間で、多くの並列プログラムでは、図4-1-3のようにノード数が増えるにつれて通信時間が増加します(そうでないプログラムもあります)。このため、ノード数を増やしていくと、並列化可能部分のCPU時間の減少よりも通信時間の増加の方がまざってしまいあるノード数で経過時間が最短となり、それ以上ノード数を増やすと逆に遅くなってしまいます。以下の例では、8台→16台では経過時間はわずかに減少していますが、16台で最短になり、それ以上はかえって遅くなっていきます。

なお、ノード数を増やして通信時間が増加する割合は、何をどのように通信するかによって異なります。プログラムによってはほとんど増加しなかったり、逆にかえって減少する場合があります。





## 4-2 どんどころが並列化できるのか

### ■ 並列性とは何か

プログラムの並列化を行う前に、そもそもMPIを使用した分散メモリ一型並列計算機(以後並列計算機)ではどのようなところが並列化できるのでしょうか？ これをまず知っておく必要があります。

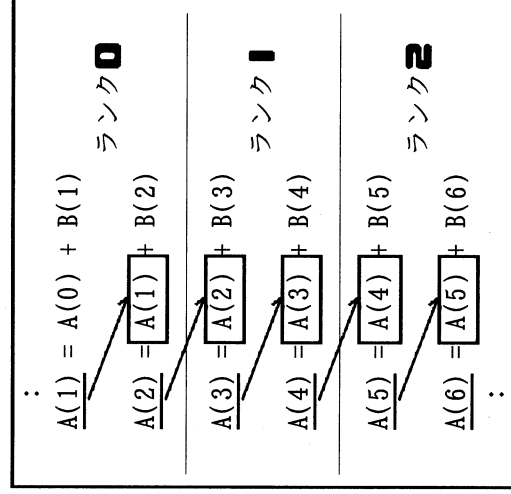
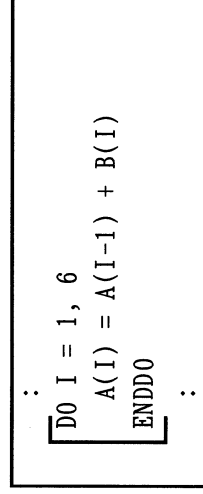
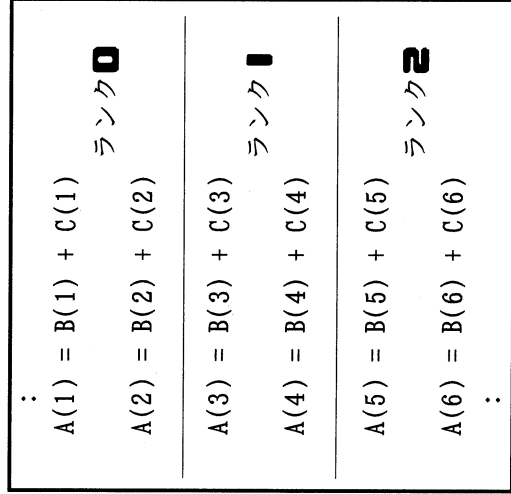
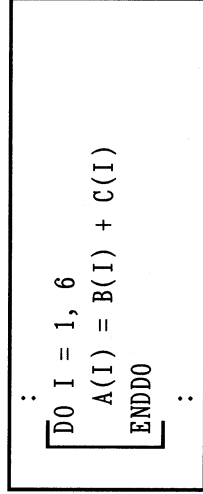
並列計算機では、並列性さえあればどんどころでも並列化することができます。

それでは並列性とは一体何なのでしょう？ 私は厳密な定義は知りませんが、D0ループの並列性について説明してみます。

D0ループの各反復間に依存性がなく独立に実行できる場合、並列性があると言います。依存関係がない場合は、ループのある反復で更新された変数(左辺の変数)が、それ以後の反復の右辺で参照されません。例えば、図4-2-1(1)のループを展開すると図4-2-1(2)のようになりますが、更新された値A(1)は以後の右辺で参照されません。従って、各ステートメントは完全に独立しており、どのような順序で実行してもかまわないわけです。このような場合は、例えば4-8(2)のように、3つのプロセスが2ステートメントずつ計算して並列化を行うことができます。

これに対し、図4-2-2(1)を展開した図4-2-2(2)では、更新されたA(1)は次の反復の右辺で参照されているため、必ず図4-2-2(2)の順序でステートメントを実行する必要があります。これを、このループには回帰参照がある、あるいは依存関係があると呼ぶこともあります。この場合、図4-2-2(2)のように並列化すると、例えばランク■のプロセスはランク□のプロセスが計算するA(2)の値が確定する前にA(2)の値を参照してしまい、結果がおかしくなってしまいます。従って、D0ループ内に回帰参照がある場合には、原則として並列化することはできません(ただし図4-2-2(1)のループは、特別な方法で並列化することができます(4-6-9節参照))。

なお、合計を求める『S = S + A(I)』や最大値を求める『X = MAX(X, A(I))』のような縮約演算は、回帰参照がありませんが並列化することができます(4-6-4節参照)。



並列計算機では、各ループに並列性があるかどうかをチェックするツールなどは提供されていませんので、ユーザーが次のような観点で自分で判断する必要があります。なお、実際のプログラムでは、並列性があるかどうかかわかりにくいD0ループはめったにありません。

- 図4-2-2(1)のように、ループ内のステートメントの左辺にA(I)がある場合、右辺にA(I-1)がないかチェックして下さい。もしある場合には回帰参照があるため、通常は並列化できません。
- ループに回帰参照があるかどうかかわからない場合、図4-2-1(2)のようにループ内のステートメントを展開して書き下したり、配列間のデータの動きを図示して調べるのも一つの方法です。
- 極めて邪道な方法ですが、単体版でループの反復順序を逆にして(例:D0 I=10,1,-1)、結果が合うかどうかを調べるという方法もあります。もしループの各反復間に依存関係があれば、ループの反復順序を逆にした場合は通常結果がおかしくなくなります。ただしこの調査方法は完全ではありませんので注意して下さい。
- D0ループの一部に回帰参照のあるステートメントがある場合、そのループ全体が並列化できなくなるのを避けるため、以下のように回帰参照のある部分だけ別のループにして、元のループを並列化するという方法もあります。

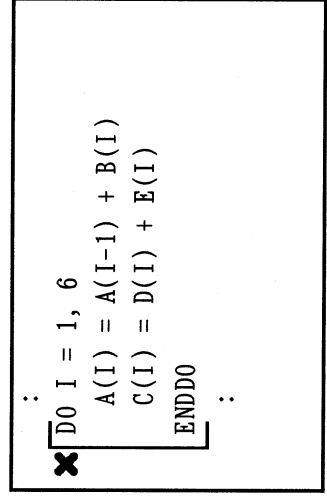


図4-2-3(1)

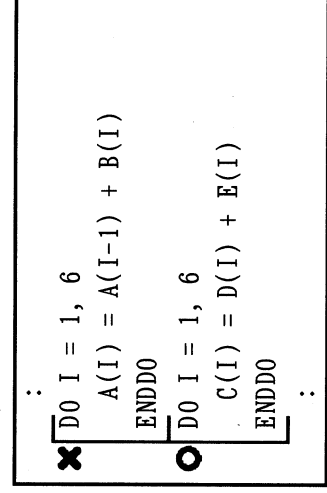


図4-2-3(2)

### ■ どんなところが並列化できるのか

本節の冒頭で、並列計算機では、並列性さえあればどんなところでも並列化することもできると述べました。その例を以下に示し、ベクトル計算機(スーパーコンピュータ)の場合と比較してみます。なお、ベクトル計算機では、コンパイラがD0ループを自動的にベクトル化します。また、ベクトル化とは並列化の一種であると覚えて下さい。

#### (1) 多重ループ

ベクトル計算機では、原則として多重ループの内側のループしかベクトル化できません(ただしコンパイラがループを入れ替えることがあります)。並列計算機では、並列化とは単に計算量を1/nに減少させることなので、図4-2-4(2)(3)に示すように、内側のループでも外側のループでも並列性さえあれば並列化が可能です。

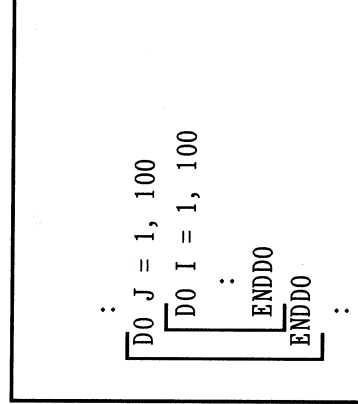


図4-2-4(1)

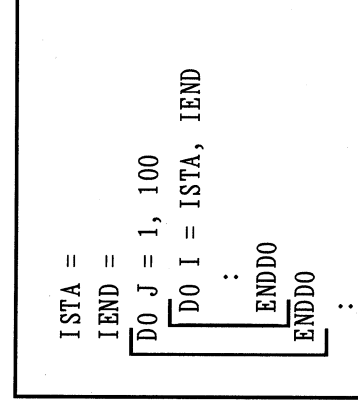


図4-2-4(2)

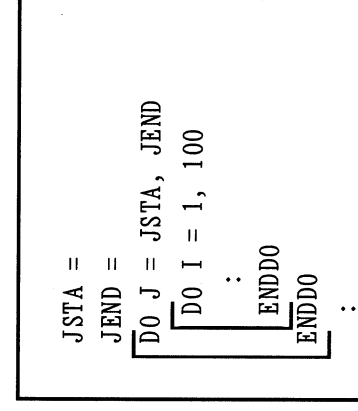


図4-2-4(3)

(2) サブルーチンコンローを含むD0ループ

ベクトルコンパイラは、各ループがベクトル化できるかどうかを解析しますが、D0ループにサブルーチンコンローを含んでいる場合、解析が複雑となるためベクトル化を行いません。一方並列計算機の場合、図4-2-5(2)のように、並列性さえあれば並列化が可能です。

```

:
DO I = 1, 100
CALL SUB
ENDDO
:
    
```

図4-2-5(1)

```

:
ISTA =
IEND =
DO I = ISTA, IEND
CALL SUB
ENDDO
:
    
```

図4-2-5(2)

(3) I/Oを含むD0ループ

ベクトル計算機の場合、I/O命令を含むD0ループはベクトル化できません(ただし、I/O命令をD0ループの外に出して、その他の演算をベクトル化することができます)。一方並列計算機の場合、図4-2-6(2)に示すように、並列性さえあれば、各プロセスが自分のローカルディスクに対してI/Oを行うことによって、I/Oそのものを並列化することも可能です。

```

:
DO I = 1, 100
WRITE (10,*)
ENDDO
:
    
```

図4-2-6(1)

```

:
ISTA =
IEND =
DO I = ISTA, IEND
WRITE (10+MYRANK,*)
ENDDO
:
    
```

図4-2-6(2)

(4) ブロック

図4-2-7(1)で処理1～処理3が独立に実行できるとします。ベクトル計算機では、D0ループの形態になっていないとベクトル化できないので、図4-2-7(1)はベクトル化できません。一方並列計算機の場合、図4-2-7(2)のようにすれば並列化は可能になります。この例として、3本の偏微分方程式(実際には連立一次方程式)を解くプログラムがあり、それぞれ独立に解くことができますが、連立一次方程式の解法自体には並列性がないとします。この場合、図4-2-7(2)のようにすれば一応並列化することができます(ただし3ノードです)。

なお、処理1～処理3が独立かどうかかわからない場合は、処理の順序を例えば 処理3→処理2→処理1 のように変えてみて、結果が合うかどうかを調べるのも(邪道ですが)1つの方法です。

```

:
処理1
処理2
処理3
:
    
```

図4-2-7(1)

```

:
IF (MYRANK == 0) THEN
  処理1
ELSEIF (MYRANK == 1) THEN
  処理2
ELSEIF (MYRANK == 2) THEN
  処理3
ENDIF
:
    
```

図4-2-7(2)

## 4-3 計算部分の並列化パターン

前節に説明したように、並列計算機の場合は並列性さえあればどんなところでも並列化は可能であり、ベクトル計算機と比較して一見いいことづくめのように見えます。しかし、どこでも並列化できるということは、裏返すと、どこを並列化するかをユーザーが決定しなければならぬということになります。

逆にベクトル計算機の場合は、基本的に多重ループの内側のD0ループしかベクトル化できなにかわりに、ほぼ自動的にコンパイラがベクトル化を行ってくれますので、ユーザーがいちいちベクトル化に関する方針を決定する必要はありません。ただし、ホットスポット(計算の集中している部分)がベクトル化されていないかかったり、ベクトル化されていても効果が出していない場合は、ベクトルチューニングを行う必要があります。

並列計算機に話を戻すと、実は、プログラムのどこを並列化するかという選択が、並列化が成功するかどうかの最大のポイントになります。この選択を誤ると、いくら頑張っても並列化による効果は得られません。

そして、これは各プログラムに固有の話になるため、ユーザーが自分で考えなくてはなりません。とはいえ、巷のプログラムも並列化という観点で言えばそうそう無限のバリエーションがあるわけではなく、マクロに見ると多くのプログラムは、これから述べるいくつかのパターンかその組み合わせで並列化できていることが分かっています。

### ■ 計算パターン1

プログラム全体のうち、極めて限られた部分がホットスポット(計算の集中している部分)になっている場合があります。例えば図4-3-1(1)のようにタイムステップのループの中に①、②、③があり、このうち②の部分が計算時間のほとんどを占め、①と③の部分はステップ数は多いのですが計算時間は短いです。また、②では配列Aの値を求めることが最終目的です。言いかえると、②で更新された配列のうち配列Aは①や③で参照または更新が行われませんが、配列Wは配列Aの計算を行うための作業配列で、②のみで使用され、①や③では参照や更新が行われません。

この場合、①と③の部分を並列化しても、ステップ数が多いため修正の手間はかかりますが、計算時間がかかっていないためパフォーマンスの向上は期待できません。このような場合、②の部分だけを並列化し、①と③の部分はブランクとして扱います。この場合の動作を以下で説明します。

● ②のD0ループの反復を図4-3-1(2)の⑤のように修正します。これによって⑤の部分が並列化、すなわち計算量が $1/n$ に縮小されました。

● ⑤で最終的に求めたいのは配列Aなので、⑤が終了した直後の⑥で、各プロセスは、配列Aのうち自分が計算した部分のデータを、他の全プロセスに送信します(図4-3-1(4)参照)。つまりお互いがお互いにデータを送り取りします。この通信はメッセージ交換ルーチンのMPI\_ALLGATHER(V)に相当します(4-6-3-2節参照)。

● ⑥の通信が終了した時点で、各プロセスは単体で実行したのと(すなわち②の終了時と)全く同じ状態になります。正確にいうと配列Wの中身は異なっていますが、前述のように配列Wは②のみで使用されないで問題ありません。

● 次に⑦、④で全プロセスが単体での③、①と全く同じ動作をします。

● 再び⑤に来ますが、⑦と④で全プロセスが単体と同じ動作をしたため、この時点で全プロセスは計算に必要な全データを持っていきます。従って、⑤に入る前にデータの送信を行う必要はありません。

このパターンは、並列化のための修正作業が少なくて済むという長所があります。一方、4-1節の最後で説明しましたが、例えば図4-3-1(4)で計算量と通信量のオーダーはどちらもnで次数が同じなので、通常は並列化すると却って遅くなってしまいます。しかし筆者の経験では、②で計算時間のほとんどを占めるようなプログラムの場合、②の中にもう1重ループが存在して計算量のオーダーが $n^2$ になっていたり、②で行う計算が非常に複雑だったりして、⑥のように通信時間のかかる派手な通信をしても並列化の効果が出る場合が多いようです。ただしノード数が増えると、次第に通信の割合が増えてサチュレートしてしまいます。

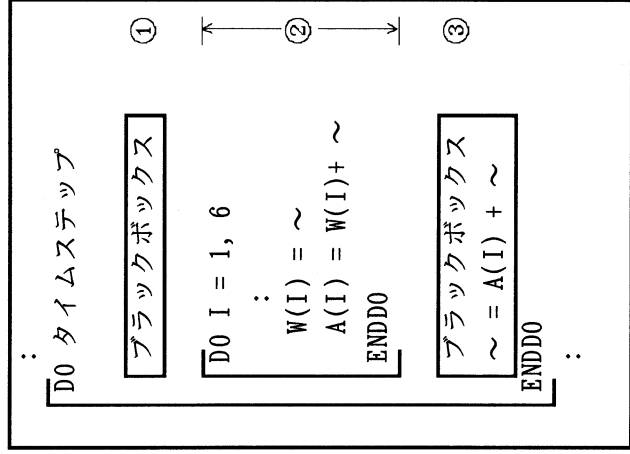


図4-3-1(1)

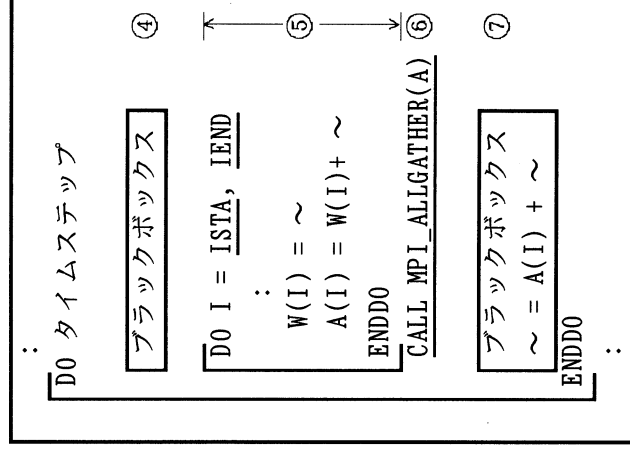


図4-3-1(2)

単体

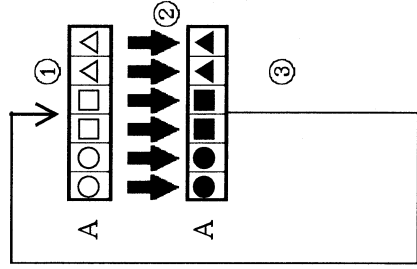


図4-3-1(3)

ランク0

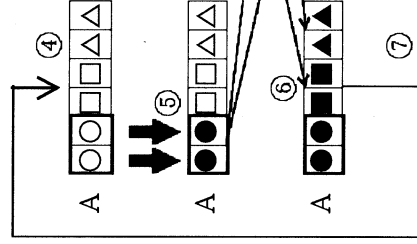
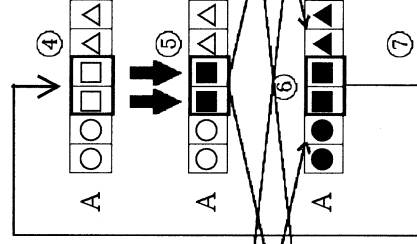


図4-3-1(4)

ランク1



ランク2

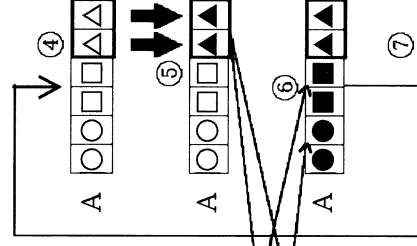


図4-3-1(4)

## ■ 計算パターン2

次のパターンは、同じような多くのD0ループがあり、ホットスポットがそれらの各ループに分散している、言いかえると顕著なホットスポットがない場合です。このようなパターンは差分法のプログラムでよく現れます。

この例を図4-3-2(1)に示します。タイムステップのループの中に①～③のループがありますが、同じようなパターンで計算量にそれほど差がありません。従って、並列化の効果をj得るためには一部のループだけでなく、タイムステップ内のほとんど全てのループを並列化する必要があります。

また、通常このパターンでは、並列化したループが終了した直後に、「計算パターン1」のような派手な通信を行うと、通信量の方がまざってしまい、並列化の効果が出ません。

以上のことから、このパターンの場合、次の方針で並列化を行います。

- (1)タイムステップのループの中には基本的に全て並列化します。
- (2)タイムステップのループが反復する間、各プロセスは自分の担当するデータのみを保持し(図4-3-2(4))、通信が必要な場合も「計算パターン1」のような通信ではなく、必要最小限の通信のみを行います。

図4-3-2を並列化した場合の動作を説明します。なお、この例は1次元の差分法(モドキ)であり、図中の『境』は固定境界、すなわちこの部分の値は変わらないことを意味します。

- ①のループの反復を④のように修正します。これによって④が並列化、すなわち計算量が $1/n$ に縮小されました。
- ②のループの反復を⑥のように修正します。⑥では配列Aを求めるのに、配列Bの自分の前後の値を使用しています。従って⑥に入る前に、自分の担当部分の境界を1つ越えたデータを隣のプロセスから送信してもらう必要があります。この通信を⑤に示すCALL SHIFT(4-6-2節と5-1節参照)で行います。  
なお、この例では、タイムステップ・ループを一周する間に配列Bの値が④で更新されているので⑤の通信が必要になりますが、配列Bが不変の値(例えば座標値など)で、タイムステップを一周しても値が更新されないのであれば、⑤の通信を行う必要はありません。
- ③のループの反復を⑦のように修正します。

このパターンは、タイムステップ・ループ内のほとんど全てのループを並列化すること、並列化に伴う通信を必要最小限しか行わないことから、通常は並列化の効果が出ます。

ところで1-1節で簡単に触れましたが、共有メモリー型並列計算機の場合、共有メモリー型並列計算機用に並列化したプログラム(コンパイラによる自動並列化、またはOpenMPの指示行による並列化:4-7-2節参照)の他に、マシン環境によってはMPIで並列化したプログラムも実行できる場合があります。差分法のようなパターン2のプログラムの場合、どちらかで並列化するのがよいか検討してみます。

パターン2のプログラムは通常各D0ループの構造が簡単なので、コンパイラによる自動並列化、またはOpenMPの指示行で、比較的簡単に並列化することができます。一方上記で説明したMPIによる並列化の場合、特にプログラムの規模が大きいと、並列化に伴う修正箇所が多くなるため、ワークロードがかかってしまいます。

従って、並列化するワークロードの観点では、共有メモリー型並列計算機用の並列化に軍配が上がりま  
す。ただしどちらがパフォーマンスが良いかは計算機の特性に依存します。

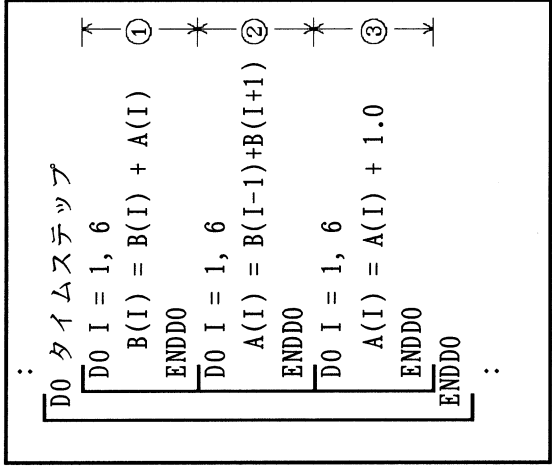


図4-3-2(1)

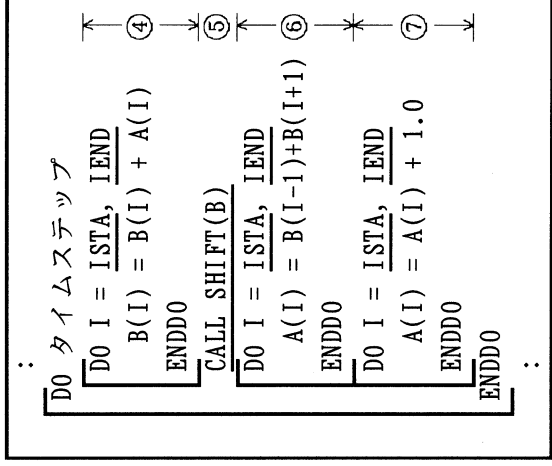


図4-3-2(2)

単体

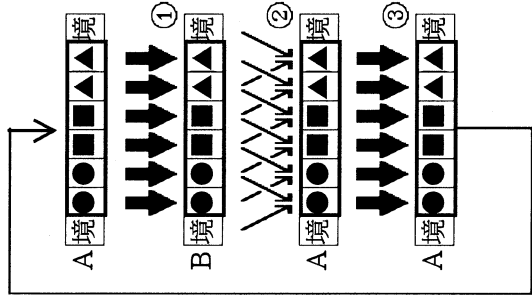
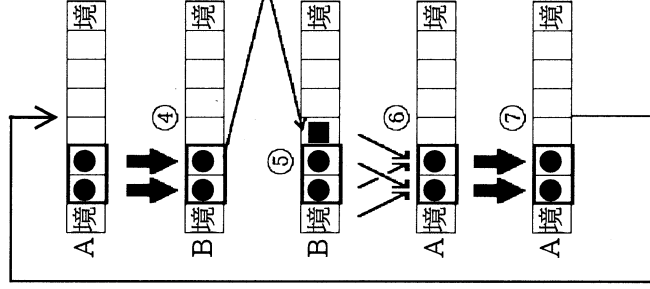
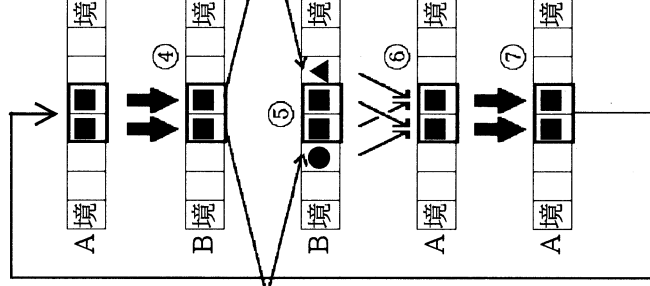


図4-3-2(3)

ランク0



ランク1



ランク2

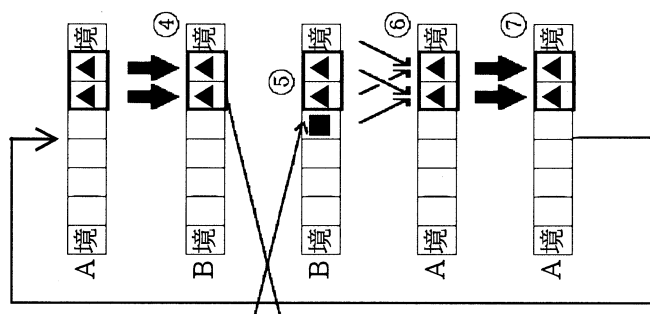


図4-3-2(4)

### ■ 計算パターン3

プログラムは一般に、図4-3-3(1)のようにメインルーチンから末端のサブルーチンまで連続と続く階層構造になっています。図4-3-3(1)の場合、■のついている末端のD0ループを図4-3-3(2)のように並列化することができず(「並」は並列化されたループ)。前述の計算パターン1や計算パターン2、あるいはベクトル計算機でのベクトル化がこれに相当します。ところで図4-3-3(1)では、上位の●の部分にも並列性があり、図4-3-3(3)のように並列化することができます。本書では、図4-3-3(2)のような並列化を末端の部分での並列化、図4-3-3(3)のような並列化を上位の部分での並列化と呼ぶことにします。なお、前者を粒度が小さい部分での並列化、後者を粒度が大きい部分での並列化と言います。

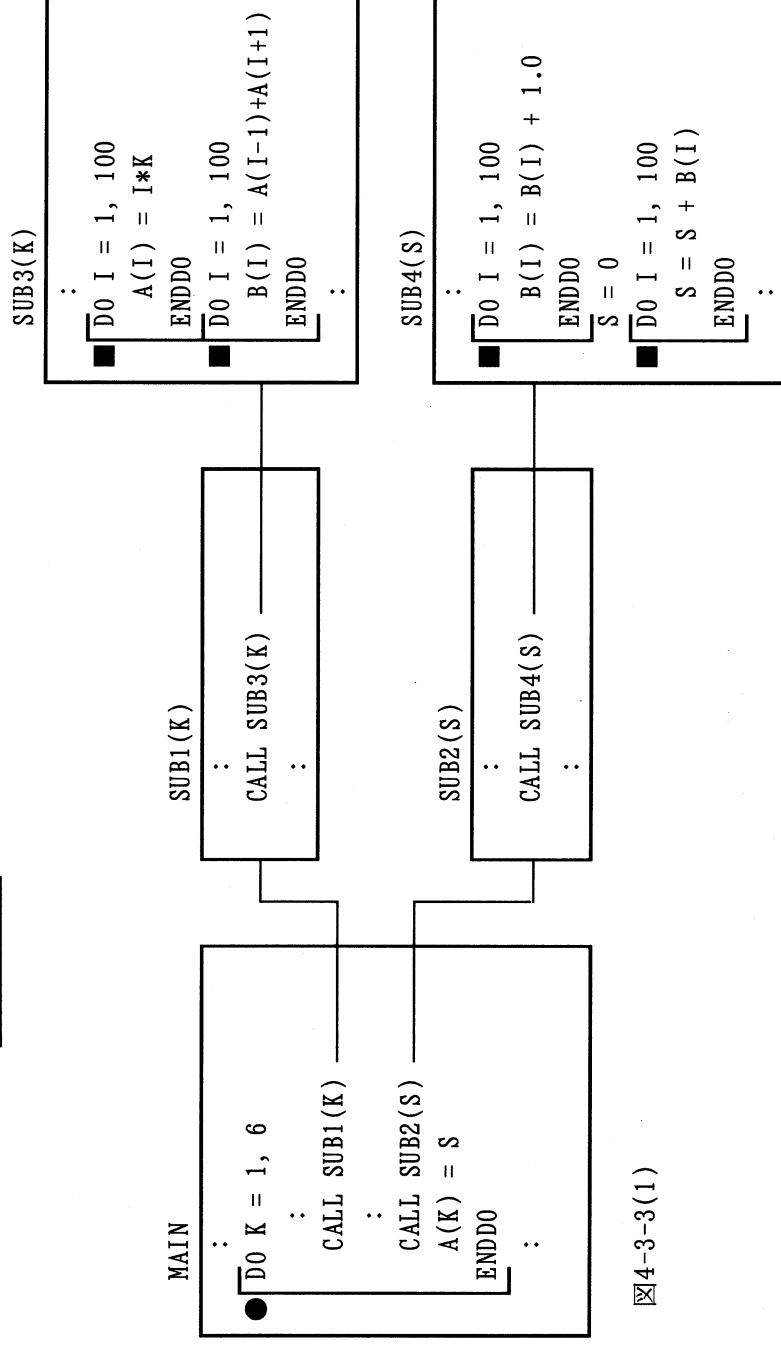


図4-3-3(1)

末端の部分のみで並列化できるプログラムもあれば、末端でも上位でも並列化できたり、あるいは末端では並列化できずに上位だけで並列化できるプログラムもあります。それでは末端でも上位でも並列化できる場合、一般にどちらで並列化するのがよいのでしょうか？ これはもちろんプログラム自体に依存しますが、図4-3-3(2)と(3)を比較すると次のことが分かります。

- 実線と二重線に示すように、上位の部分で並列化した場合は並列化のための修正量は少ないですが、末端の部分で並列化した場合、各ループを並列化しなければならぬため、修正量が多くなります。
- 二重線に示すように、上位の部分で並列化した場合は計算の最後に1回のみ通信を行います、末端の部分で並列化した場合は何度も通信する必要があり、通信の立上り時間がかかります。

以上のことから、どこを並列化するかを検討する際、末端の部分で並列化できる場合でも、上位の部分で並列化できないかどうか必ず検討して下さい。上位のループに並列性があるかどうか分からない場合、4-2節で述べた邪道な方法として、ループの反復順序を逆にして調べるという手もあります。なお、タイムステップや収束反復のループには並列性はありません。

『上位の部分と末端の部分の両方で並列化することはできないのか?』という質問が出る場合があります、上位の部分で並列化すれば末端では並列化する必要はなく、その逆も同様です。

上位の部分で並列化できるプログラムとして、次のような例があります。

- 連立一次方程式を例えればガウスの消去法で解いているプログラムで、解くべき連立一次方程式自体が100本あり、各連立一次方程式が独立に解けるとします。この場合、1回のガウスの消去法自体を並列化するのが「末端の部分の並列化」にあたります。一方例えば4プロセスの場合、各プロセスに25本ずつ連立一次方程式を割り当てると「上位の部分の並列化」になります。



● 例えば流体の計算を行うプログラムで、境界の流速を(時速1m, 2m, ...)と変えて計算し、それぞれの条件で得られた物理量を最後にまとめるようなプログラムがあります。通常は前述の『計算パターン2』で末端の部分の並列化を行います。各流速での計算を独立に行うことができ、上位の部分で並列化を行うことができます。もっともこの例などは、例えば単体ジョブを4本同時に実行し、各ジョブは100個の境界の流速のデータのうち25個のデータのみを処理するのとほとんど同じであり、並列ジョブとは言えないかもしれませんが...

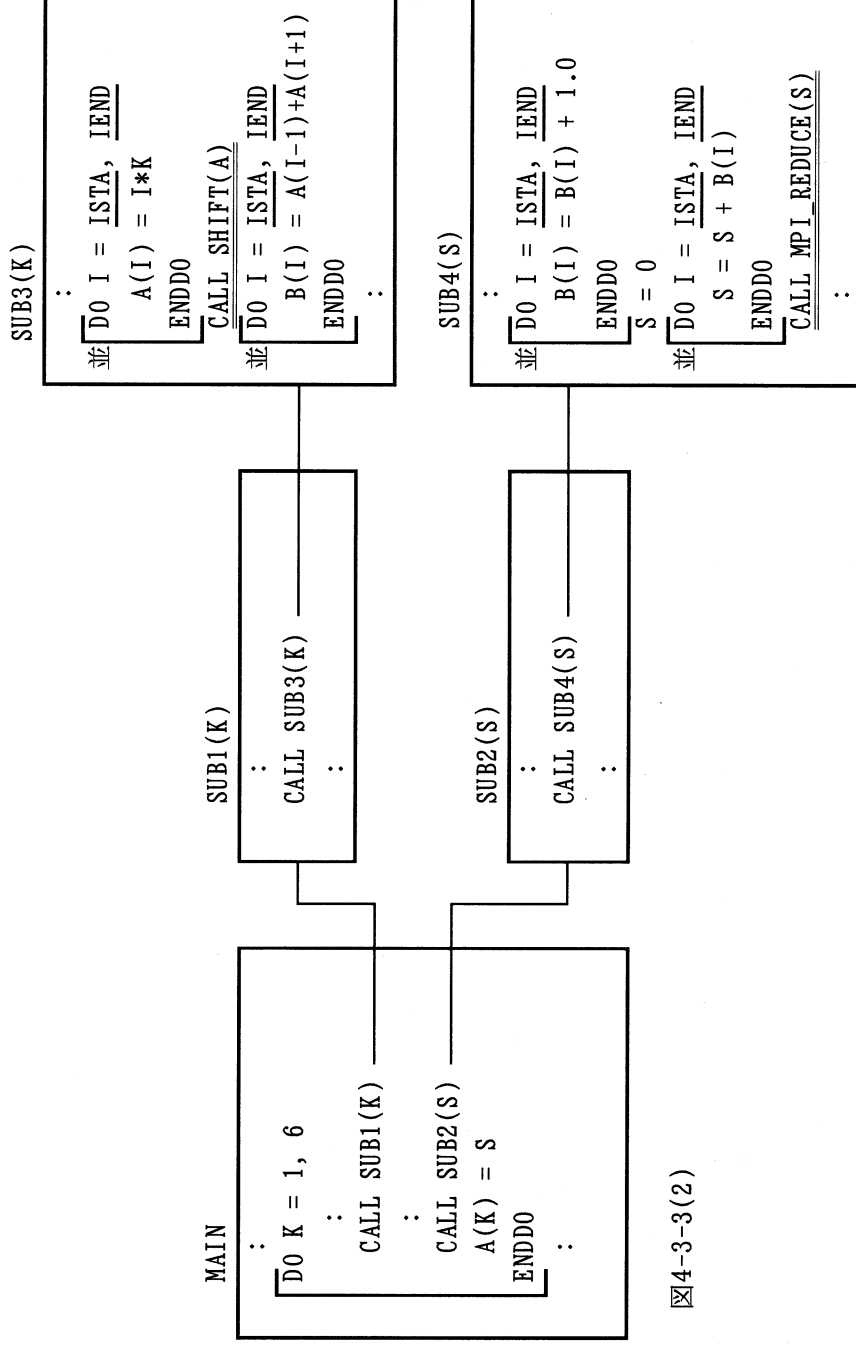


図4-3-3(2)

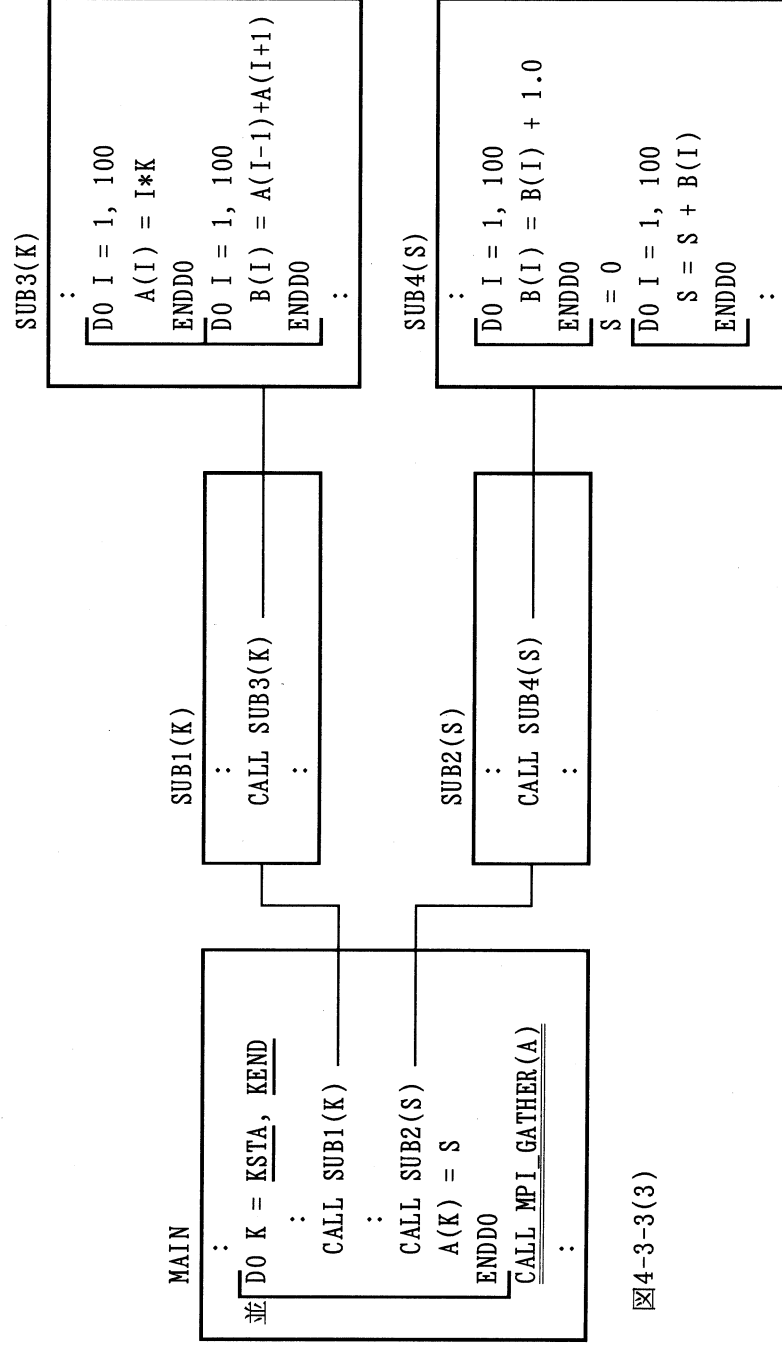


図4-3-3(3)

## 4-4 並列化に伴う入出力部分の修正

計算部分の並列化に伴い、入出力部分の修正が必要になります。本節では修正方法について説明します。なお、ファイルは順次編成ファイルを想定します。

### 4-4-1 並列化に伴う入力部分の修正

まず入力部分(READ文)の修正方法について説明します。

#### ■ 標準入力

MPIで並列化したプログラムの場合、標準入力(装置番号5番)から読み込みを行うと、マシン環境によってはエラーになる場合があります。標準入力以外のお勧めできません。標準入力以外の変更に、以下に示す他の入力パターンで読み込みを行うようにして下さい。なおマシン環境によっては、標準入力から読み込むための環境変数などが提供されている場合がありますが、並列化したプログラムを他のマシン環境へ移植する可能性がある場合は、やはり標準入力以外の変更に良いでしょう。

#### ■ 入力パターン1

図4-4-1(1)に示すように、「0 1 2」の入った入力ファイルが共有ディスクに置かれていて、各プロセスが「0 1 2」の全てのデータを使用したい場合(★の付いた場合(★のみを使用したい場合も同様))、図4-4-1(3)のREAD文(標準入力以外)を何も修正せずに実行すると、各プロセスが「0 1 2」を読み込みます。図中の⑩は装置番号を示します。

図4-4-1(1)(3)の方法はプログラムを修正する必要がないので簡単ですが、以下の欠点があります。

- 入力ファイルが共有ディスクにあるので、読み込む際にプロセス間で競合が発生する。
- 読み込んだデータ共有ディスクのファイルサーバーから各ノードに通信するための通信時間がかかる。

一方図4-4-1(2)では、あらかじめ、入力ファイルをrpcコマンドなどで各ノードのローカルディスクにコピーし、各プロセスは自分のローカルディスクからデータを読み込むので、読み込み時間が短縮されます。大量の入力データを読み込まなければならなかったり、プロセス数が非常に多くて読み込みの競合が異常に発生して入力データを読み込む時間が無視できないほど長い場合は、図4-4-1(2)の方がよいかもかもしれません。ただしファイルをあらかじめローカルディスクにコピーする手間がかかります。

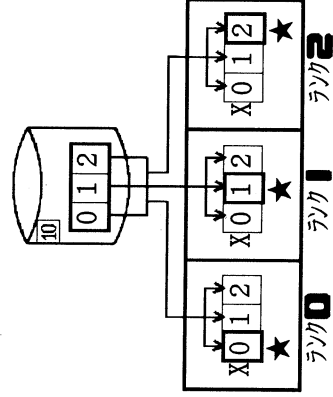


図4-4-1(1)

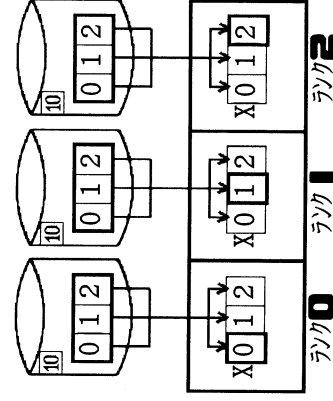


図4-4-1(2)

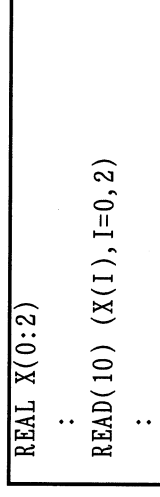


図4-4-1(3)

■ 入力パターンの2

図4-4-2(1)(2)に示すように、例えば**ランク0**のプロセスが代表して入力データを読み込み、その全てをMPI\_BCASTで他のプロセスに送信する(細い実線部分)という方法もあります。この方法は通信のオーバーヘッドが発生し、またプログラムによっては、たくさんの配列を送る必要があるため修正が面倒になるという欠点があります。なお、市販の並列プログラムでは、この方式になっているプログラムが多いようです。各プロセスが全入力データのうちの**一部しか使用しないのであれば**、図4-4-2(3)(4)のように、例えば**ランク0**のプロセスが、読み込んだデータの**うち各プロセスが必要とする部分のみを**、MPI\_SCATTERなどを使用して送信する(細い実線部分)という方法もあります。なお、MPI\_SCATTERを使用した場合、配列Xとは別に配列Yが必要となりますが、配列Yを使用しない方法もあります(3-3-5節、3-8-1節、4-6-3節参照)。

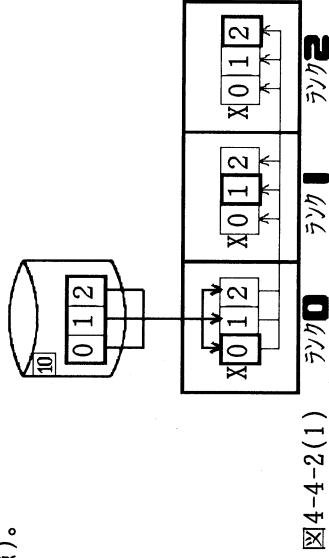


図4-4-2(1)

```
REAL X(0:2)
:
IF (MYRANK==0) READ(10) (X(I), I=0,2)
CALL MPI_BCAST(X,3,MPI_REAL,
& 0,MPI_COMM_WORLD,IERR)
&
```

図4-4-2(2)

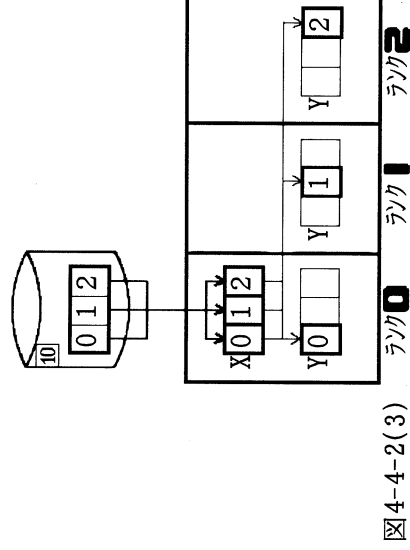


図4-4-2(3)

```
REAL X(0:2),Y(0:2)
:
IF (MYRANK==0) READ(10) (X(I), I=0,2)
CALL MPI_SCATTER(X,1,MPI_REAL,
& Y(MYRANK),1,MPI_REAL,
& 0,MPI_COMM_WORLD,IERR)
```

図4-4-2(4)

図4-4-2(5)に示すように、入力データを読み込んで、エラーを発見したらエラーメッセージを表示して処理を終了する単体プログラムを並列化し、図4-4-2(6)のように**ランク0**が代表して入力データを読み込むようにした場合、**ランク0**はエラーを発見し、STOP文を終了しますが、他のプロセスは**①**で**ランク0**から送られるデータの受信待ちになり、終了しません。全プロセスを終了させたい場合は、図4-4-2(7)に示すようにMPI\_ABORT(付録参照)を使用します。「入力パターン1」の場合は、図4-4-2(8)に示すように全プロセスが入力データを読み込んでエラーを発見し、STOP文を終了するので、MPI\_ABORTを実行する必要はありません。

```
READ(10,*) A
IF (エラー) THEN
PRINT *, '==ERROR=='
STOP
ENDIF
:
```

図4-4-2(5) 単体プログラム

```
:
IF (MYRANK==0) THEN
READ(10,*) A
IF (エラー) THEN
PRINT *, '==ERROR=='
CALL MPI_FINALIZE(~)
STOP
ENDIF
ENDIF
CALL MPI_BCAST(A,~)
:
```

図4-4-2(6) ✕

```
:
IF (MYRANK==0) THEN
READ(10,*) A
IF (エラー) THEN
PRINT *, '==ERROR=='
CALL MPI_ABORT
& (MPI_COMM_WORLD,9,IERR)
CALL MPI_FINALIZE(~)
STOP
ENDIF
ENDIF
CALL MPI_BCAST(A,~)
:
```

図4-4-2(7) ○

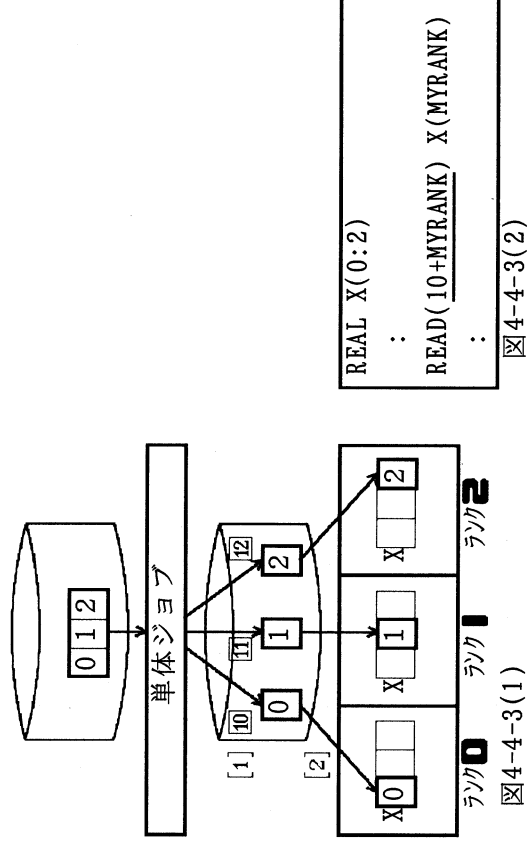
```
READ(10,*) A
IF (エラー) THEN
IF (MYRANK==0)
PRINT *, '==ERROR=='
CALL MPI_FINALIZE(~)
STOP
ENDIF
:
```

図4-4-2(8) ○

### ■ 入力パターン3

並列ジョブを実行する前に、図4-4-3(1)の[1]に示すように、単体ジョブを実行し、入力ファイルを並列ジョブで各プロセスが担当するデータごとに別ファイルに分けます。分割したファイルを共有ディスクに作成する場合は、下線に示すようにファイルの装置番号をランクの値で区別します。各プロセスのローカルディスクに作成する場合は同一の装置番号で構いません。

次に[2]と図4-4-3(2)に示すように、並列ジョブを実行し、各プロセスは自分の担当するデータの入ったファイルから読み込みを行います。なお、ファイルをNFSマウントされた共有ディスクに作成する場合は、上記の単体ジョブと並列ジョブを連続に行います。[1]で書き出したデータが、[2]を開始した時点ではキャッシュに入っているファイルに入っておらず、[2]の読み込みで誤動作する可能性がありますので、ファイルにデータが入ったのを確認してから並列ジョブを実行するようにして下さい。

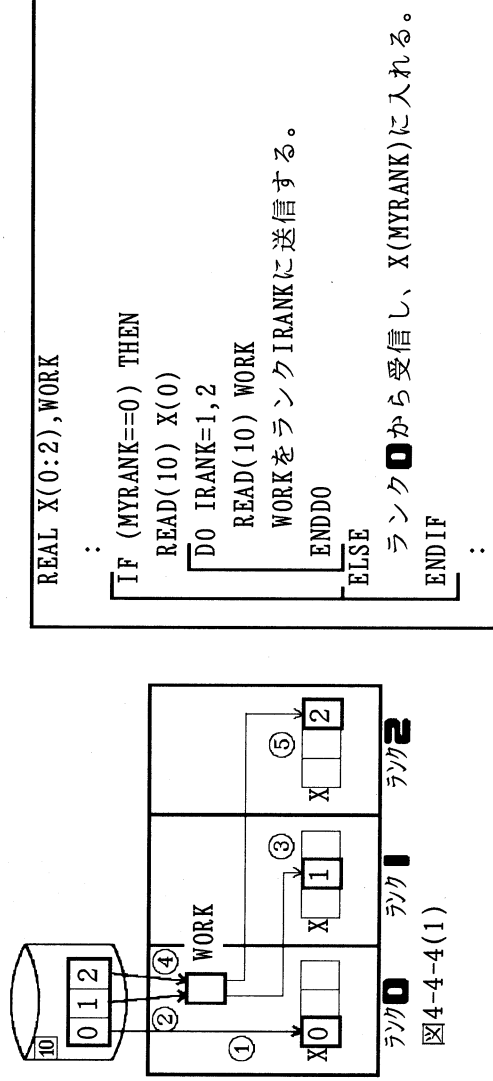


```
REAL X(0:2)
:
READ(10+MYRANK) X(MYRANK)
:
```

図4-4-3(2)

### ■ 入力パターン4

図4-4-4(1)で、まずランク0が①で自分の担当範囲をファイルから読み込みます。次にランク0は②でランク1の担当範囲を読み込んでランク1に送信し、ランク1は③で受信します。最後にランク0は④でランク2の担当範囲を読み込んでランク2に送信し、ランク2は⑤で受信します。プログラムの例を図4-4-4(2)に示します。



```
REAL X(0:2), WORK
:
IF (MYRANK==0) THEN
  READ(10) X(0)
  DO IRANK=1,2
    READ(10) WORK
    WORKをランクIRANKに送信する。
  ENDDO
ELSE
  ランク0から受信し、X(MYRANK)に入れる。
ENDIF
:
```

図4-4-4(2)

■ 入力パターン5

図4-4-5(1)で、各プロセスは共有ディスクに作成された入力データのうち、自分が担当する部分のみを読み込み、それ以外の部分は空読みします。バイナリーファイルの場合、以下の(2)と(3)の方法で作成されたファイルは、それぞれ(4)と(5)の方法で読み込みます(空読みする部分を下線で示します)。

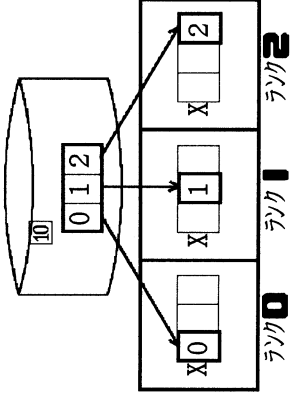


図4-4-5(1)

```

REAL X(0:2)
:
WRITE(10) (X(I), I=0, 2)
または
WRITE(10) X
:
    
```

図4-4-5(2)

```

REAL X(0:2)
:
READ(10) (DUMMY, I=0, MYRANK-1),
          (X(I), I=MYRANK, MYRANK),
          (DUMMY, I=MYRANK+1, 2)
:
    
```

← 空読み  
← 空読み

図4-4-5(4)

```

REAL X(0:2)
:
DO I=0, 2
WRITE(10) X(I)
ENDDO
:
    
```

図4-4-5(3)

```

REAL X(0:2)
:
DO I=0, MYRANK-1
READ(10) ← 空読み
ENDDO
DO I=MYRANK, MYRANK
READ(10) X(I)
ENDDO
DO I=MYRANK+1, 2
READ(10) ← 空読み
ENDDO
:
    
```

図4-4-5(5)

■ 入力パターン6

3-1節で紹介したMPIの拡張機能であるMPI-2では、MPI-I0(7イ・オー)というI/Oルーチン群が提供されています。この機能が使用できるかどうかはマシン環境に依存します。本書はMPI-I0の範囲で説明しているので、MPI-I0の説明は省略します。

## 4-4-2 並列化に伴う出力部分の修正

次に出力部分(WRITE文, PRINT文)の修正方法について説明します。

### ■ 標準出力

図4-4-6(1)のような標準出力(WRITE(6)またはPRINT文)は、並列に実行すると全プロセスが同じ内容を早い者順に出力します。そこで例えばランク0のみが出力するように、図4-4-6(2)の下線部を追加します。並列化のデバッグの際、図4-4-6(3)のようにわざと全プロセスに標準出力を書かせることがあります。この場合、出力行数が多かったり各行の文字数が多いと、各行が混ざって表示が乱れることがあります。

プログラムによっては、多くの箇所で標準出力を行っている、修正が非常に面倒になる場合があります。マシン環境によっては、図4-4-6(1)のまま、以下のように書き出すことのできる環境変数などが提供されています。これはMPI自体の機能ではなく、マシン環境(並列処理環境)が提供する機能です。

- 標準出力を、ランクごとに別のファイルに書き出す。

- 指定したランクのプロセスのみが標準出力を書き出す。

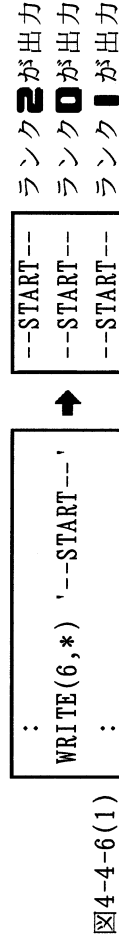


図4-4-6(1)

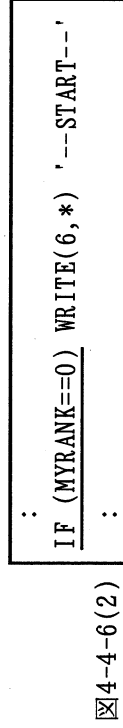


図4-4-6(2)

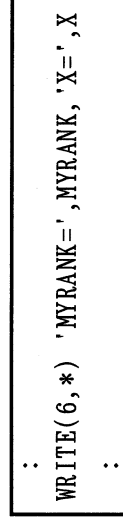


図4-4-6(3)

### ■ 出力パターン1

標準出力以外の出力ファイル(順次編成ファイルを想定)に関しては、図4-4-7(1)(2)のように複数のプロセスが共有ディスク上の同じファイルに単に書き出しを行うと、中身がグチャグチャになってしまいます。

この場合、図4-4-8(1)(2)のように、各プロセスは自分が計算した部分のデータを、MPI\_GATHER(V)を使用して例えばランク0のプロセスに送り(細い実線部分)、ランク0のプロセスが代表して出力ファイルに書き出します。なお、MPI\_GATHERを使用した場合に配列Yが必要となりますが、配列Yを使わずランク0の配列Xを集める方法もあります(4-6-3-1節参照)。

また、配列Xを縮小(4-5-7節)している場合の出力方法については4-5-7-3節で説明します。

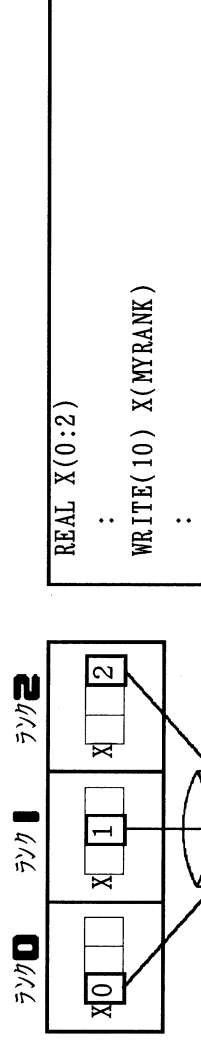


図4-4-7(1)

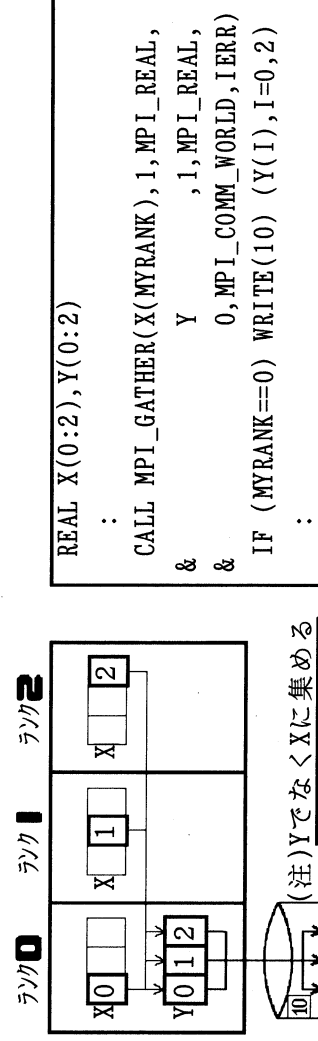


図4-4-8(1)

■ ランク0が代表して出力する方法で誤動作する例

出力パターン1のようにランク0のプロセスが代表して出力を行う場合、誤動作する例を紹介합니다。まず並列化する前の単体プログラムを図4-4-9(1)に示します。実行前に10番ファイルが存在していない場合、①が失敗し、err=3に従って③に飛び、④～⑥を実行して10番ファイルを作成します。

このプログラムを並列化し、ランク0が代表して出力するように⑦と⑧を追加したのが図4-4-9(2)です。実行前に10番ファイルが存在しない場合、「全プロセスが①を失敗し、err=3に従って③に飛び、ランク0だけが④～⑥を実行して10番ファイルを作成する。」のが正常な動作です。ところが、たまたまランク0が他のプロセス(例えばランク1)より先行して④～⑥で10番ファイルを作成し、その後でランク1が①を実行したとすると、ランク1は①が成功して②を実行し、上記の正常な動作と異なってしまいます。

この場合、図4-4-9(3)の図に示すように、同期をとるだけの集団通信ルーチンMPI\_BARRIER(付録参照)を挿入します。すると、実行前に10番ファイルが存在していない場合、ランク0がランク1より先行しても、ランク0は⑧で(ランク1を含む)全プロセスが⑧に到着するまで待つので、後から①に到達したランク1は①を失敗して②を実行せずに③に飛び、正常に動作します。

```

:
OPEN(10, STATUS='OLD', err=3) ①
計算 ②
3 CONTINUE ③
OPEN(10, STATUS='NEW') ④
WRITE(10,*) Y ⑤
CLOSE(10) ⑥
:

```

図4-4-9(1)

```

:
OPEN(10, STATUS='OLD', err=3) ①
計算 ②
3 CONTINUE ③
IF (MYRANK==0) THEN ⑦
OPEN(10, STATUS='NEW') ④
WRITE(10,*) Y ⑤
CLOSE(10) ⑥
ENDIF ⑧
:

```

図4-4-9(2)

```

:
OPEN(10, STATUS='OLD', err=3) ①
計算 ②
3 CONTINUE ③
CALL MPI_BARRIER(~) ⑨
IF (MYRANK==0) THEN ⑦
OPEN(10, STATUS='NEW') ④
WRITE(10,*) Y ⑤
CLOSE(10) ⑥
ENDIF ⑧
:

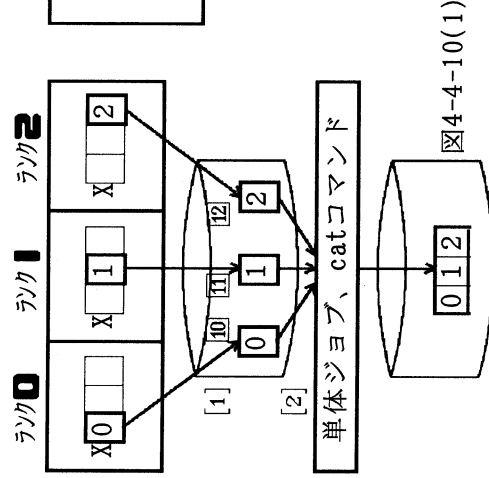
```

図4-4-9(3)

■ 出力パターン2

この方法は入力パターン3の反対です。図4-4-10(1)の[1]と図4-4-10(2)に示すように、並列ジョブを実行し、各プロセスは自分が計算した部分のデータをファイルに書き出します。ファイルは共有ディスクに作成する場合は、下線に示すようにファイルの装置番号をランクの値で区別します。各プロセスのローカルディスクに作成する場合は同一の装置番号で構いません。

次に[2]に示すように、各プロセスが作成したファイルを、単体ジョブやcatコマンドで1つのファイルに合体します。なお、ファイルをNFSマウントされた共有ディスクに作成する場合、上記の並列ジョブと単体ジョブを連続に実行すると、[1]で書き出したデータが、[2]を開始した時点ではキャッシュに入っていてファイルに入っておらず、[2]の読み込みで誤動作する可能性がありますので、ファイルにデータが入ったのを確認してから単体ジョブを実行するようにして下さい。



### ■ 出力パターン3

この方法は入力パターン4の反対です。図4-4-11(1)で、まずランク0が①で自分の担当範囲をファイルに書き出します。次にランク1が②で自分の担当範囲をランク0に送信し、ランク0は③でファイルに書き出します。最後にランク2が④で自分の担当範囲をランク0に送信し、ランク0は⑤でファイルに書き出します。プログラム例を図4-4-11(2)に示します。詳細は4-5-7-3節を参照して下さい。

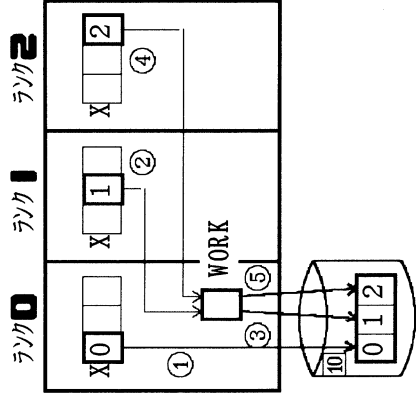


図4-4-11(1)

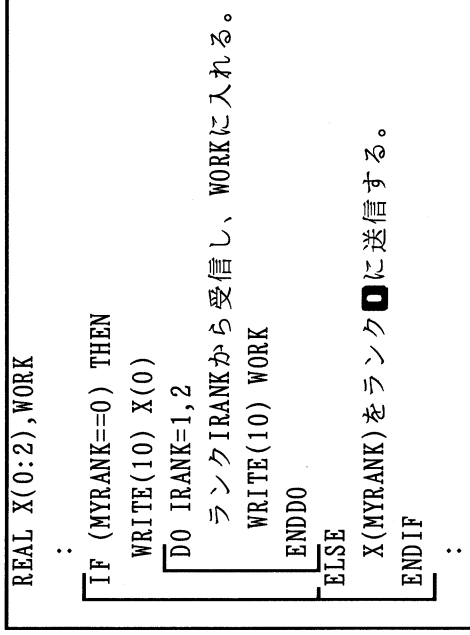


図4-4-11(2)

### ■ 出力パターン4

入力パターン6で説明したように、MPI-2では、MPI-IO(71・0-)というI/Oルーチン群が提供されています。この機能が使用できるかどうかはマシン環境に依存します。本書はMPI-1の範囲で説明しているので、MPI-10の説明は省略します。



### 4-4-3 入出力に関するその他の考慮点

● I/Oが非常に多いプログラムの場合、図4-4-12に示すように、CPU時間をいくら並列化して速くしてもジョブ全体としてはパフォーマンスは向上しません。このような場合、もしI/O自体に並列性があるのであれば、4-2節で述べたように、各プロセスが自分のローカルディスクにI/Oを行うことによって、I/O自体を並列化することができます。

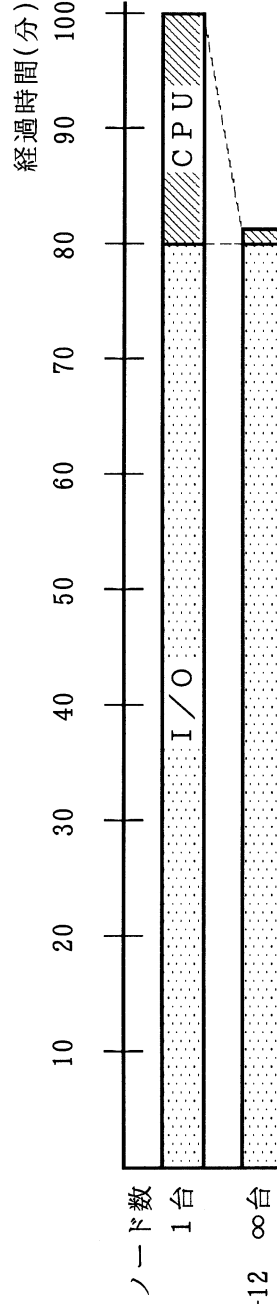


図4-4-12

● 単一のマシンではメモリが足りないため、本来ならメモリー上に展開できるデータを作業ファイルを使用して計算しているためにI/Oが多くなっている場合があります。この場合、4-5-7節で説明するように、並列化の際に配列を縮小すれば、1つの並列ジョブが全ノードのメモリーを使用できるようになります。それによってメモリーが足りるようになれば、作業ファイルを使用しないようにプログラムを修正してI/Oを減らすことができます。

● 一般にプログラムは図4-4-13のように3つの部分に分けられます。このうち入出力部分は経過時間がほとんどかからないので並列化せず、計算部分のみを並列化するのが一般的です。

4-4-2節の「出力パターン1」で、WRITE文をランク0のプロセスだけが書き出すようにIF文を付加する方法を説明しました。図4-4-13のようにWRITE文が多くていちいちIF文を付加するのが面倒な場合は、①で必要なデータをランク0のプロセスに送信した後、③でランク0以外のプロセスが終了するようになれば、追加するIF文は③の1行で済みます。なお、③を実行するよりも前に、②のMPI\_FINALIZEを必ず1度実行しなければなりません。②を実行せずに③を実行した場合、マシンの環境によっては誤動作しますので注意して下さい(「付録」のMPI\_FINALIZE参照)。

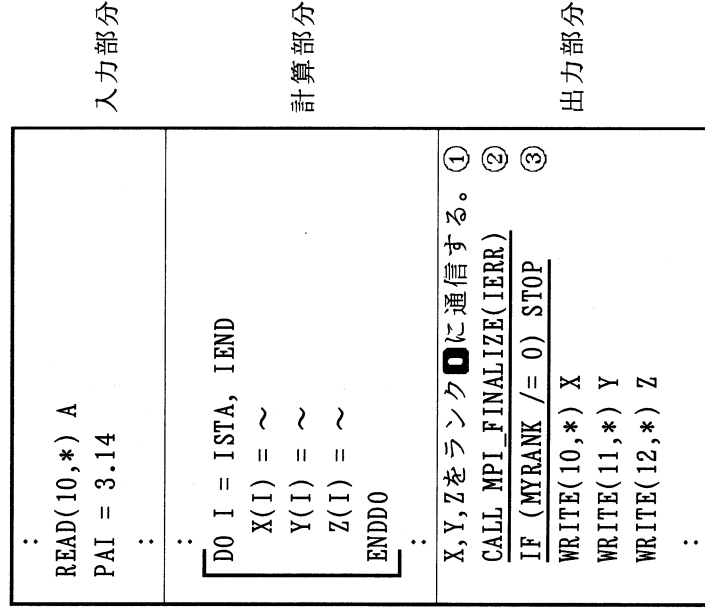


図4-4-13

## 4-5 ループの分割方法

### 4-5-1 ループ反復と配列

4-2節で、『並列計算機では、並列性さえあればどこでも並列化できる』と説明しました。とはいえ一般的にはほとんどの場合、D0ループを並列化します。D0ループの各反復で計算量が同じであれば、各プロセスの仕事を均等にするためには、(全反復回数)÷(プロセス数)(回)の反復を各プロセスに割り当てればよいことになります。例えば図4-5-1(1)(3)では、ループの全反復回数が9回でプロセス数が3つなので、各プロセスには  $9 \div 3 = 3$ (回)の反復を割り当てています。このうち図4-5-1(1)は、D0ループの反復とループ内で使用している配列A、B、Cが1対1に対応しています(図4-5-1(2)参照)。これに対し図4-5-1(3)では、指標ベクトルNを使用しているため、D0ループの反復とループ内で使用している配列Aは1対1に対応していません。このような場合、計算を各プロセスに割り当てる方法として、図4-5-1(4)の右図のようにD0ループの反復を各プロセスに分ける方法のほかに、図4-5-1(5)のように配列Aで分割する方法もあります。本節では、D0ループの反復を各プロセスに分ける方法を説明します。図4-5-1(1)のように、ループの反復と配列が1対1に対応している方が、分割方法を図示できて説明がしやすいのでこちらを想定します。従って『D0ループの反復を分ける』のと『配列を分割する』のは同じ意味で使用します。なお、図4-5-1(3)(4)を並列化する方法については4-6-6節で説明します。

単体

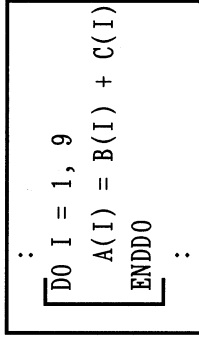


図4-5-1(1)

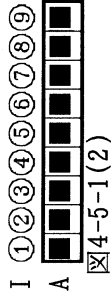


図4-5-1(2)

単体

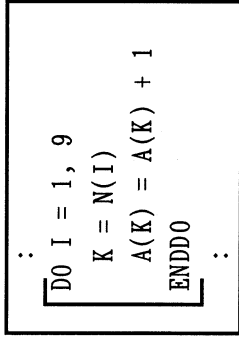


図4-5-1(3)

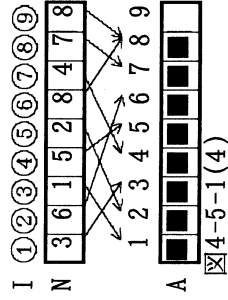
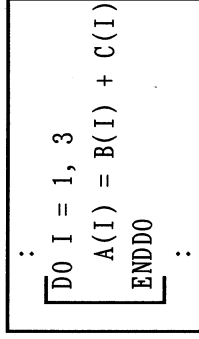
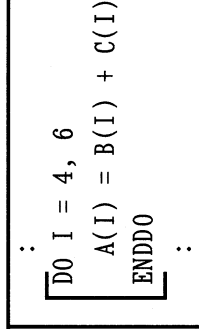


図4-5-1(4)

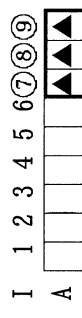
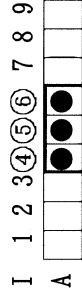
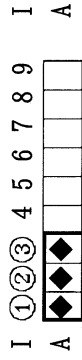
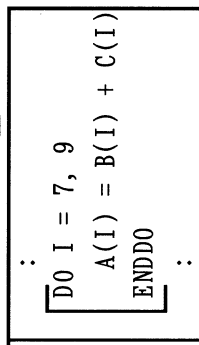
ランク0



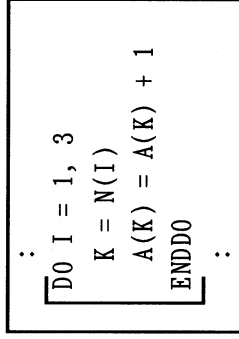
ランク1



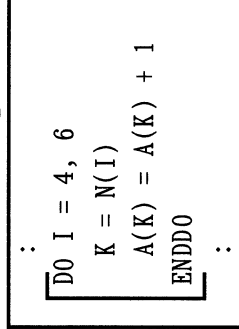
ランク2



ランク0



ランク1



ランク2

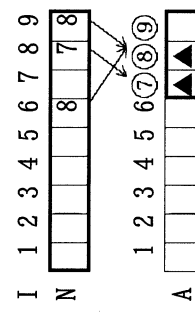
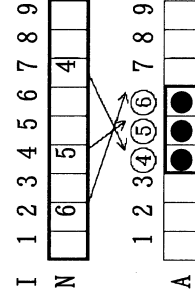
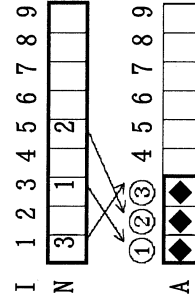
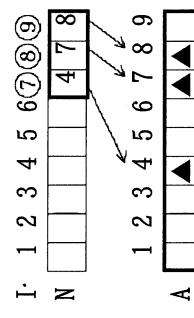
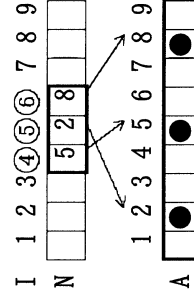
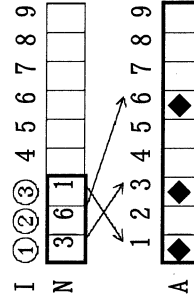
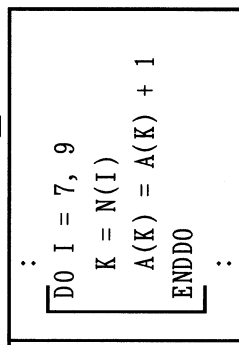


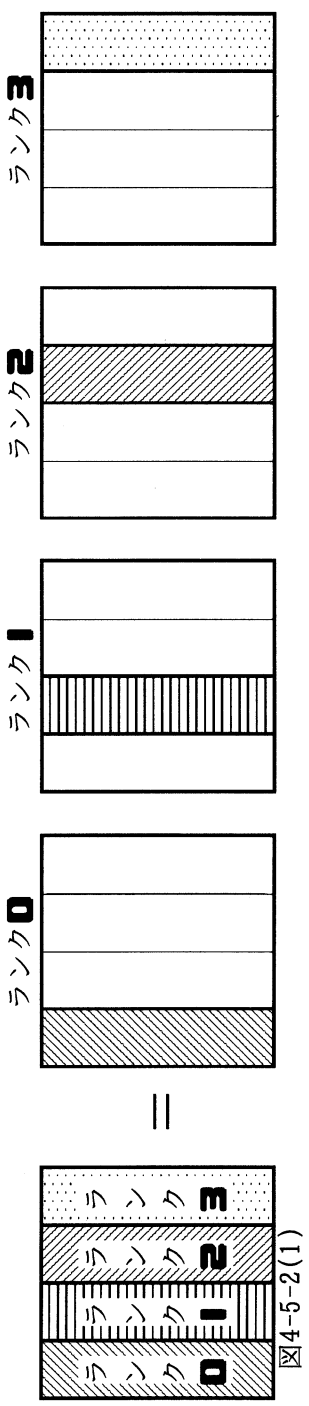
図4-5-1(5)

## 4-5-2 配列の分割方法

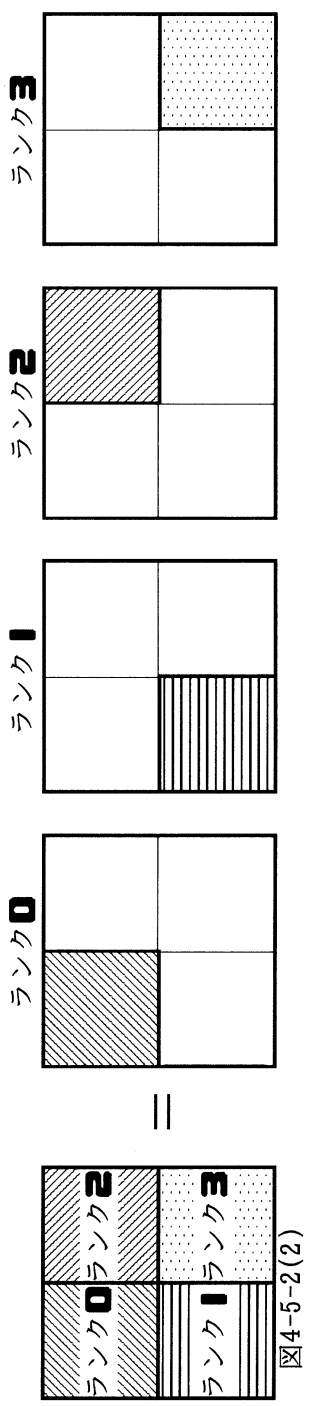
2次元配列を想定し、配列の分割方法について説明します。

### ■ ブロック分割

例えばプロセス数が4の場合、各プロセスに1/4ずつの連続した領域を割り当てる方法をブロック分割と呼びます。大半の並列化はブロック分割で行います。なお、2次元配列の場合、図4-5-2(1)のように2次元目(列方向)を分割する方法と、1次元目(行方向)を分割する方法があり、並列化によって発生する通信量を考慮して分割の方向を決定します(4-5-8節参照)。

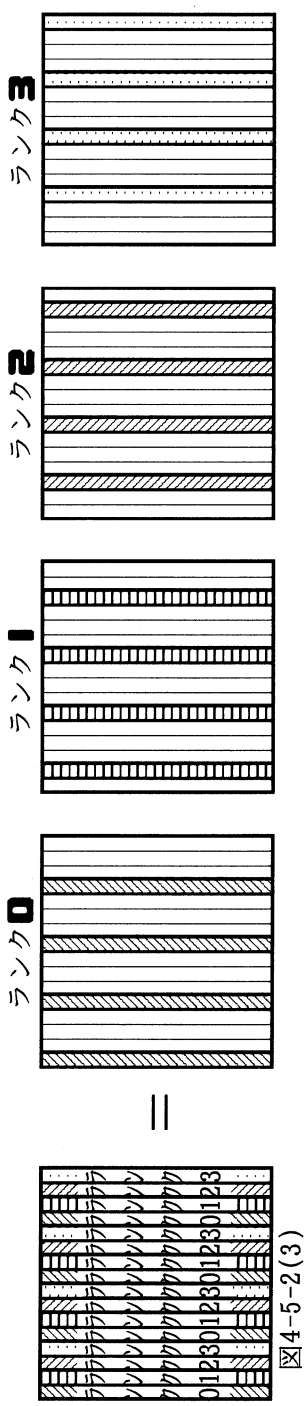


2次元配列の場合、1方向(行方向または列方向)のみでブロック分割するほかに、図4-5-2(2)に示すように両方向でブロック分割する方法もあります。この方法は、差分法などで、1方向のみのブロック分割よりも通信量を減らしたい場合に使用します(4-5-8節参照)。



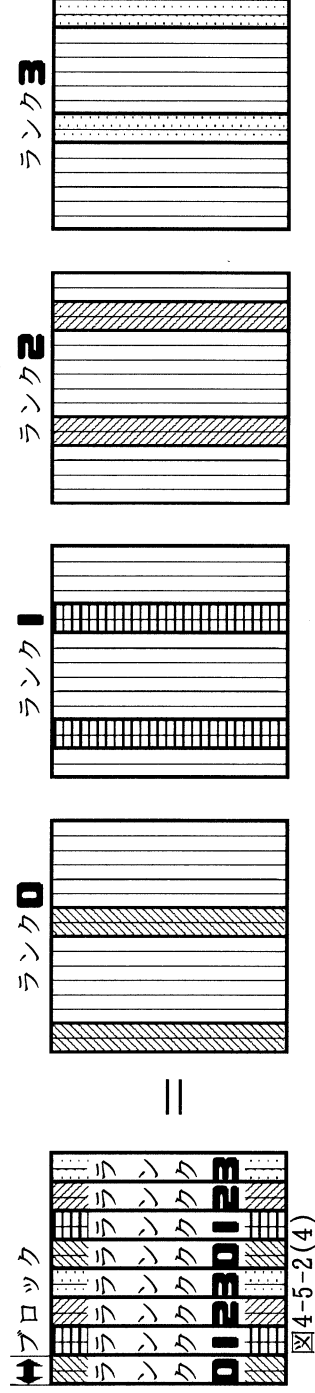
### ■ サイクリック分割

図4-5-2(3)のように、各列(または行)を、ランク0,1,2,3,0,1,2,3,...のように各プロセスに割り当てる方法をサイクリック分割と呼びます。ブロック分割では何らかの理由でプロセス間のロードバラン



## ■ ブロック・サイクリック分割

図4-5-2(4)に示すように、何列(または行)かをまとめて1つのブロックとし、そのブロックを単位としてサイクリック分割を行う方法を、ブロック・サイクリック分割と呼びます。全ての列をブロックとした場合はブロック分割に、1列のみをブロックとした場合はサイクリック分割になりますので、ブロック分割とサイクリック分割はブロック・サイクリック分割の両極端であると考えることもできます。

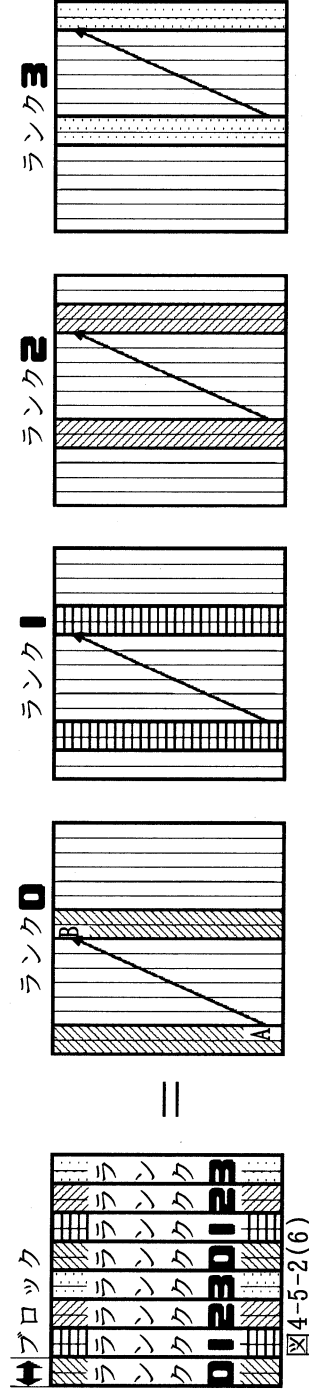
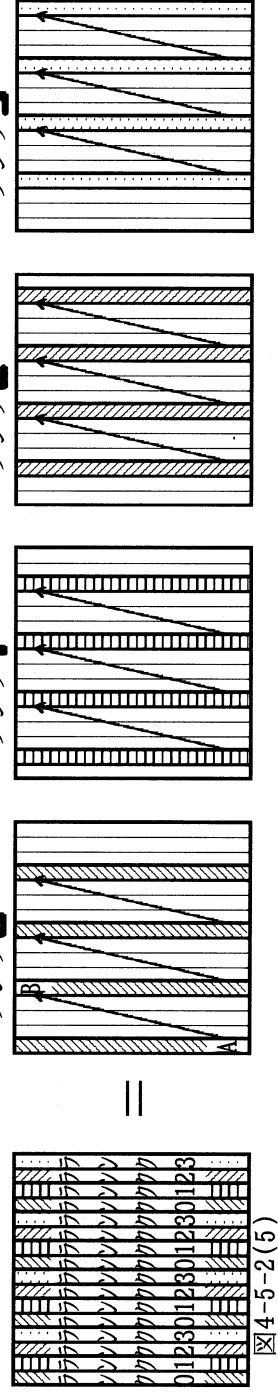


ブロック・サイクリック分割は、「本当はサイクリック分割にしたいが、サイクリック分割にすると別の問題が発生するため、多少サイクリック性を犠牲にする代わりにその問題を改善する」場合に使用します。以下に例を示します。

### (1) キャッシュミスの減少

図4-5-2(5)に示すように、配列を縮小(4-5-7節参照)せずにサイクリック分割を行い、要素Aの次に要素Bを処理するとします。この場合、AとBのメモリー内の距離が離れているため、キャッシュ・ミス(参考文献[6]参照)が発生してパフォーマンスが低下する可能性があります。キャッシュミスは図中の↗の部分で発生します。

ブロック・サイクリック分割にした図4-5-2(6)では、↗の数が減るため、キャッシュミスの頻度が減ります。



### (2) 通信量の減少

図4-5-2(5)(6)で、隣のプロセスとの境界(太い縦線部分)で横方向の通信が発生するとします。ブロック・サイクリック分割の方が境界の数が少ないので、通信量と通信回数が共に減少します。

### 4-5-3 ブロック分割でロードバランスが不均等になる例

本節では、ブロック分割では各プロセスのロードバランスが不均等になる例を紹介いたします。

#### ■ 例 1

大気の汚染を解析するプログラムがあり、汚染の激しい所ほど計算量が多いとします。このプログラムで東京湾の周辺を計算したところ、例えば図4-5-3(2)のようになっています。黒い部分は汚染が激しい(つまり計算量が多い)ことを意味します。図4-5-3(2)の計算量を列ごとに合計すると、図4-5-3(1)のようになります。

この場合、図4-5-3(2)の上段のように、列方向にブロック分割すると、図4-5-3(3)のようにプロセス間でのロードバランスが著しく不均等になります。ところがサイクリック分割にすると、図4-5-3(3)のようにロードバランスはプロセス間でほぼ均等になります。

通常、自然現象はなめらかに変化するので、個々の列では計算量がバラついていても、サイクリック分割にして合計すれば、(よほど運が悪くない限り)プロセス間でのバラつきは相殺されます。

ただし、図4-5-3(2)で横方向の通信が必要なプログラムの場合、サイクリック分割にすると他のプロセスとの境界が増えるので、通信量が多くなってしまいます(前節で説明したブロック・サイクリック分割にすれば、境界の数が減るので通信量のある程度減らすことができます)。

各列の計算量

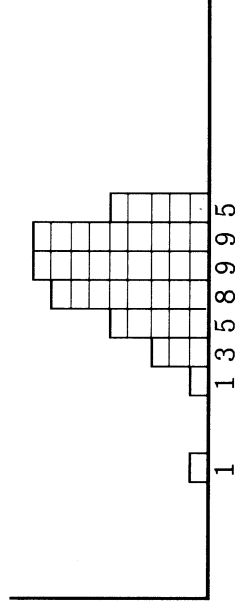


図4-5-3(1)

ブロック分割  
サイクリック分割

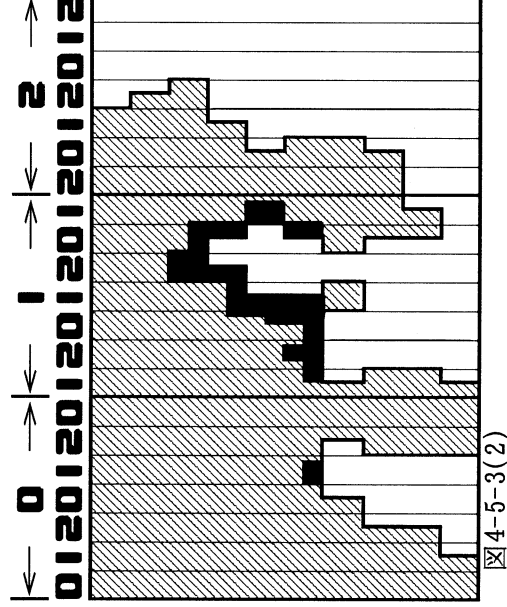


図4-5-3(2)

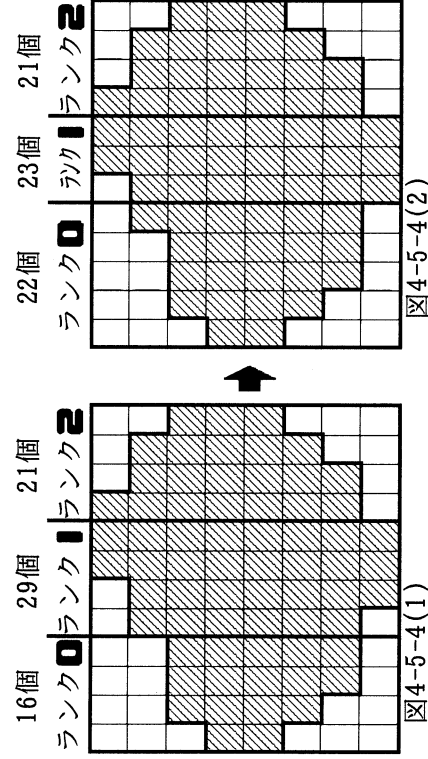
	ランク0	ランク1	ランク2
ブロック分割	1	40	0
サイクリック分割	14	15	12

図4-5-3(3) 計算量の合計

■ 例 2

図4-5-4(1)の2次元領域のうち、実際に計算するのは着色した格子部分のみだとします(例えばば湖)。この場合、例えば図4-5-4(1)のように列方向にブロック分割すると、ロードバランスが不均等になります。また解法が差分法の場合、一般に横方向の通信が必要になるので、例1のようにサイクリック分割にすると、他のプロセスとの境界が増えて通信量が多くなってしまいます。

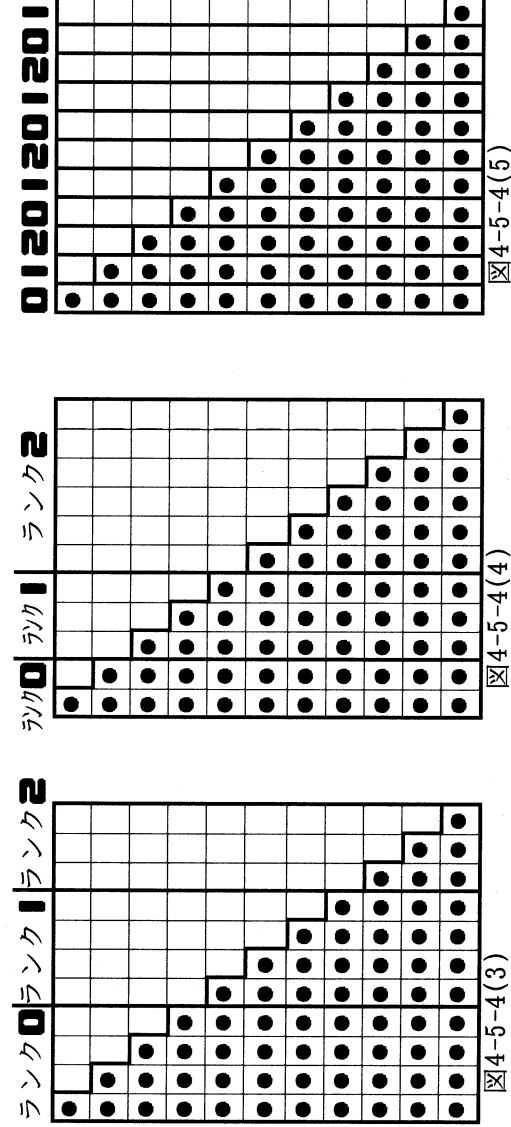
この場合、図4-5-4(2)のように、実際に計算する格子の数が各プロセスでほぼ均等になるようにブロック分割するという方法が考えられます。



■ 例 3

図4-5-4(3)に示すように、●の要素だけを計算する三角行列の場合、列方向にブロック分割するとロードバランスは不均等になります。

この場合、図4-5-4(4)のように、各プロセスが担当する要素数がほぼ均等になるようにブロック分割するか、または図4-5-4(5)のようにサイクリック分割(またはブロック分割)にします。



■ 例 4

図4-5-5(1)は例えば1000×1000の行列を表します。実際に計算する部分(太線内)が、最初は1000×1000(正方行列全体)、次が999×999、次が998×998、...のように、計算が進むにつれて変化するとします(例えば5-3節の密行列の連立一次方程式のLU分解)。

これを列方向にブロック分割すると、最初は各プロセスのロードバランスは均等ですが、計算が進むにつれてランク0の計算部分が終了し、さらにランク1の計算部分も終了し、不均等になってしまいます。

これを図4-5-5(2)のようにサイクリック分割(またはブロック・サイクリック分割)にすれば、計算が進んでもロードバランスは常にほぼ一定になります。

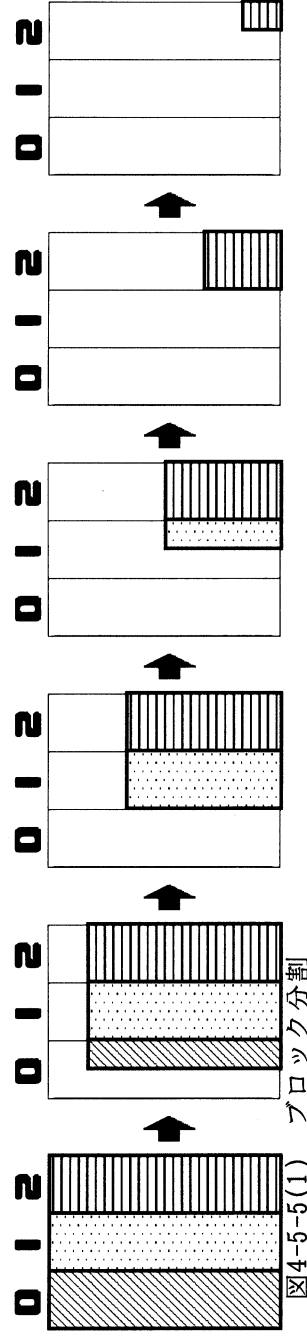


図4-5-5(1) ブロック分割

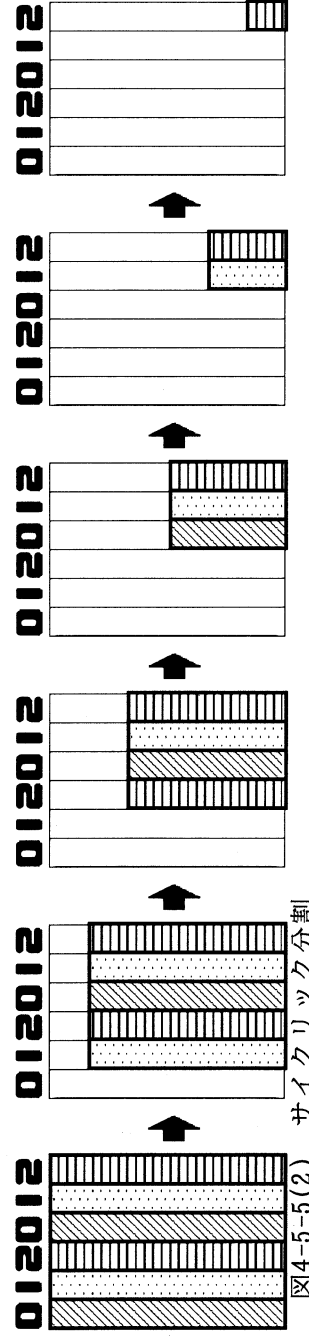


図4-5-5(2) サイクリック分割

■ 例 5

箱の中に入った粒子の動きを調べるプログラムを考えます。各粒子は、自分を中心とした円内(点線)の粒子のみと力が働き、力、速度、位置の計算を行います。例えば図4-5-4(1)の粒子(16)は、円内に他の粒子が4個あるため多くの計算を行います。一方粒子(1)は、円内に他の粒子がないため計算をほとんど行いません。

また、粒子の存在する位置が箱の中央に近いほど、円内にたくさん粒子があるとなります。

この例では、箱の左端の粒子から順に粒子の番号が付いています。粒子番号でブロック分割した場合、図4-5-6(1)に示すように、箱の端付近の粒子を担当するランク0と2のプロセスは計算量が少なく、中央付近の粒子を担当するランク1のプロセスは計算量が多くなり、ロードバランスが不均等になってしまいます。この場合、図4-5-6(2)に示すように、粒子番号でサイクリック分割(またはブロック・サイクリック分割)にすれば、各プロセスが担当する粒子の(箱内の)位置が偏らなくなるため、ロードバランスはほぼ均等になります。

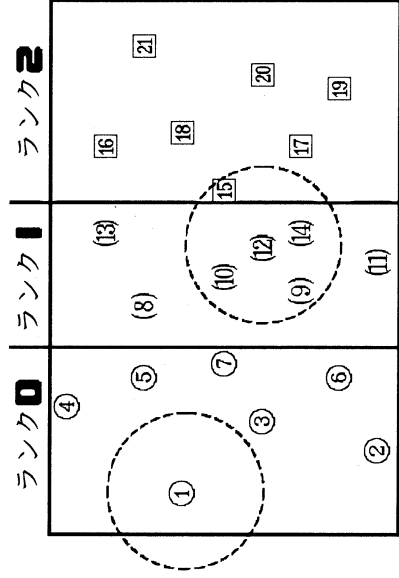


図4-5-6(1)

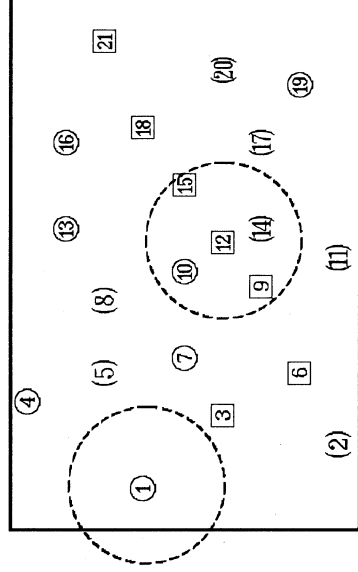


図4-5-6(2)

○:ランク0  
( ):ランク1  
□:ランク2

## 4-5-4 ブロック分割

D0ループを各プロセスに分割する場合、いままでは説明を簡単にするため、1~9までの9個の要素を3つのプロセスに分割するというように、具体的な値がわかっていて、しかもちようど割りきれられる例で説明しました(図4-5-7(1))。ところが実際のプログラムでは、図4-5-7(2)のように、D0ループの反復が未知数になっていたり、プロセス数を3個でも4個でも分割できるようにしたいという場合もあります。本節では、処理する要素数、プロセス数がともに任意の値で、要素数がプロセス数で割りきれない可能性がある場合のブロック分割方法について説明します。

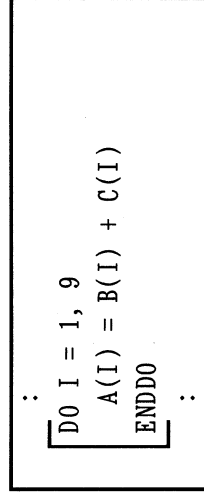


図4-5-7(1)

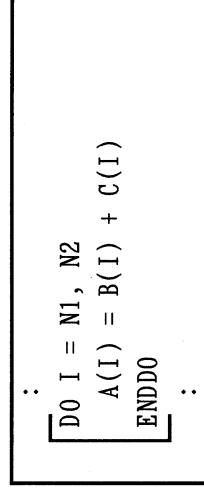


図4-5-7(2)

### ■ 分割方法の種類

ブロック分割には2つの方法があります。これを要素が①~⑥までの6個、プロセス数が4個(ランク0~3)の場合で説明します。この例では要素数がプロセス数で割りきれないことに注意して下さい。

#### (1) 分割方法1

最初の分割方法では、要素を各プロセスに図4-5-7(3)のように分配します。この意味ですが、各プロセスが図4-5-7(4)のようにそれぞれ箱をもってしているとします。そして要素を1つずつ矢印の順序で箱に入れていき、要素を全部箱に入れ終わった状態が図4-5-7(3)ということとなります。

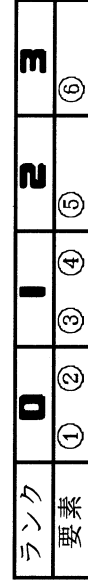


図4-5-7(3)

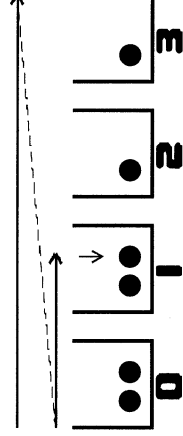


図4-5-7(4)

#### (2) 分割方法2

2つ目の分割方法では、図4-5-7(5)に示すように、要素数が全プロセス数で割りきれない場合、最後以外のプロセスは同じ要素数を担当し、最後のプロセスだけがそれより少ない要素を担当します。

なお、図4-5-7(6)のように要素数がプロセス数より少しだけ多い場合、最後以外のプロセス(ランク2)も、少ない要素数になることがあります。

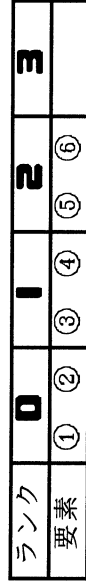


図4-5-7(5) 6個の要素を4プロセスで分割

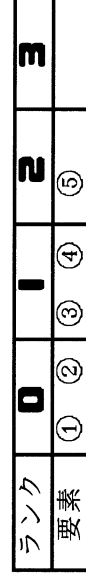


図4-5-7(6) 5個の要素を4プロセスで分割

### ■ 分割方法1と2の比較

分割方法1と2のどちらを選択するかは好みによります。並列ジョブでは最も処理時間の長いプロセスに足をひっぱられますが、図4-5-7(3)(5)から分かるように、どちらの分割方法でも各プロセスが担当する最多の要素数は同じなので、パフォーマンスはほとんど同じです。



## ■ 実際の分割方法

実際のプログラムでは、分割を行うためのサブルーチンを使用するのが便利です。筆者が作成したサブルーチン例を以下に示します。例えば、要素番号の最小値 $n1(=1)$ から最大値 $n2(=20)$ までを、プロセス数 $nprocs(=4)$ で分割したとき、ランクの値が $irank(=0)$ のプロセスの担当範囲の下限 $ista(=1)$ と上限 $iend(=5)$ が戻ります(割りきれない場合も考慮されます)。

分割方法1と分割方法2に対応するサブルーチンをそれぞれ図4-5-8(1)と図4-5-8(2)に示します。

```
CALL PARA_RANGE(n1, n2, nprocs, irank, ista, iend)
(値を指定する) 1 ↘ 20 ↘ 4 ↘ 0 ↘ 1 ↘ 5 (値が戻る)
```

- $n1$  : 整数。分割したい要素番号の最小値を指定します。
- $n2$  : 整数。分割したい要素番号の最大値を指定します。
- $nprocs$  : 整数。分割に使用する全プロセス数を指定します。
- $irank$  : 整数。担当範囲を調べたいプロセスのランクを指定します。 $0 \leq irank \leq (nprocs-1)$ です。
- $ista$  : 整数。ランク $irank$ のプロセスが担当する要素の下限が戻ります。
- $iend$  : 整数。ランク $irank$ のプロセスが担当する要素の上限が戻ります。

```
SUBROUTINE PARA_RANGE(N1,N2,NPROCS,IRANK,ISTA,IEND)
  IWORK1 = (N2-N1+1)/NPROCS
  IWORK2 = MOD(N2-N1+1,NPROCS)
  ISTA = IRANK*IWORK1 + N1 + MIN(IRANK, IWORK2)
  IEND = ISTA + IWORK1 - 1
  IF (IWORK2 > IRANK) IEND = IEND + 1
END
```

図4-5-8(1) 分割方法1に対応するサブルーチン

```
SUBROUTINE PARA_RANGE(N1,N2,NPROCS,IRANK,ISTA,IEND)
  IWORK = (N2-N1)/NPROCS + 1
  ISTA = MIN(IRANK*IWORK + N1, N2+1)
  IEND = MIN(ISTA+IWORK-1, N2)
END
```

図4-5-8(2) 分割方法2に対応するサブルーチン

なお、1つのプロセスが担当する要素数( $IEND-ISTA+1$ )が少なすぎると、並列プログラムによっては誤動作します(4-8-2節の「プログラムの並列化方法のミス(2)」参照)。各プロセスが担当する要素数が決定した時点で、正常に動作する最低限の要素数以上かどうかをチェックするようにして下さい(4-8-1節の「(5) 並列に動作させる最低限の修正」参照)。

## ■ 実際の適用例

図4-5-9(1)のA(1)~A(N)(N=1000)までの合計を求めるプログラムをサブルーチン**PARA\_RANGE**(分割方法1)を使用して並列化した例を図4-5-9(2)に、実行結果を図4-5-9(3)に示します。

実行結果から分かるように、1000個の要素を3プロセスでブロック分割した場合、ピタリは割りきれませんがほぼ均等に分割が行われ、またプロセス数がいくつでも同じ計算結果となっています。

```
PROGRAM MAIN
PARAMETER (N=1000)
DIMENSION A(N)
DO I = 1, N
  A(I) = I
ENDDO
SUM = 0.0
DO I = 1, N
  SUM = SUM + A(I)
ENDDO
PRINT *, 'SUM = ', SUM
END
```

図4-5-9(1)

```
PROGRAM MAIN
INCLUDE 'mpif.h'
PARAMETER (N=1000)
DIMENSION A(N)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD,NPROCS,IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD,MYRANK,IERR)
CALL PARA_RANGE
& (1,N,NPROCS,MYRANK,ISTA,IEND)
DO I = ISTA, IEND
  A(I) = I
ENDDO
SUM = 0.0
DO I = ISTA, IEND
  SUM = SUM + A(I)
ENDDO
CALL MPI_REDUCE(SUM,SSUM,1,MPI_REAL,
& MPI_SUM,0,MPI_COMM_WORLD,IERR)
SUM = SSUM
IF (MYRANK == 0) PRINT *, 'SUM = ', SUM
PRINT *, 'MYRANK = ', MYRANK,
& ' ISTA = ', ISTA, ' IEND = ', IEND
CALL MPI_FINALIZE(IERR)
END
```

図4-5-9(2)

```
[aoyama@node01]/u/aoyama: mpirun -np 2 a.out
SUM = 500500.000
MYRANK = 0 ISTA = 1 IEND = 500
MYRANK = 1 ISTA = 501 IEND = 1000
[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
SUM = 500500.000
MYRANK = 0 ISTA = 1 IEND = 334
MYRANK = 1 ISTA = 335 IEND = 667
MYRANK = 2 ISTA = 668 IEND = 1000
```

図4-5-9(3)

補足ですが、例えば図4-5-9(4)の①のループを分割してISTAとIENDを分割して②のループに対しては図4-5-9(5)の下線部のように分割します(図4-5-9(6)参照)。

```
DO I=1,100 ①
  A(I) = ~
ENDDO
DO I=30,90 ②
  A(I) = ~
ENDDO
```

図4-5-9(4)

```
DO I=ISTA,IEND
  A(I) = ~
ENDDO
DO I=MAX(ISTA,30),MIN(IEND,90)
  A(I) = ~
ENDDO
```

図4-5-9(5)

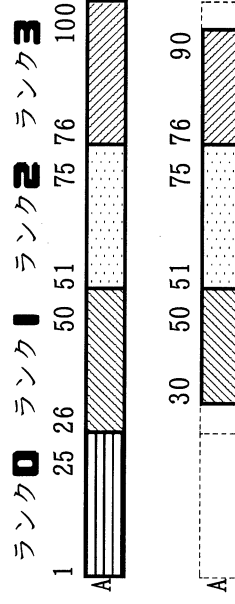


図4-5-9(6)

## 4-5-5 サイクリック分割

4-5-3節で説明したように、ブロック分割ではプロセス間のロードバランスが不均等になり、サイクリック分割では均等になる場合があります。本節では、サイクリック分割の具体的な方法を説明します。

### ■ 分割方法1

図4-5-10(1)のD0ループを、プロセス数NPROCSでサイクリック分割する方法を図4-5-10(2)に、3プロセスでの分割の様子を図4-5-10(3)に示します。ループ反復の下限が「MYRANK+1」なので、各プロセスは下限の1つずつずれた位置から間隔NPROCSで反復します。この方法は最も簡単なので、特に問題がなければこの方法を用います。

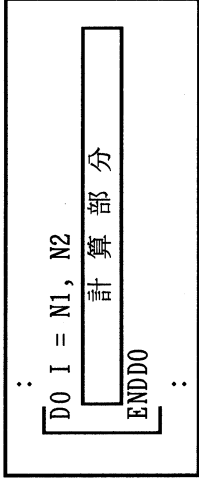


図4-5-10(1) 元のプログラム

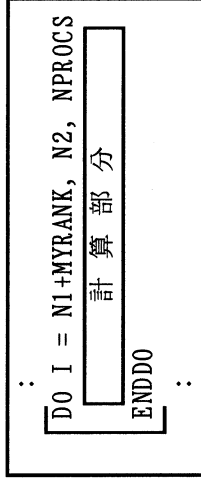


図4-5-10(2)

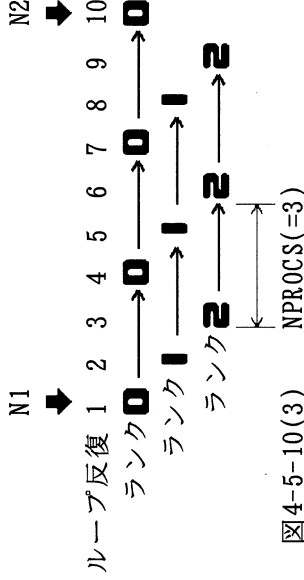


図4-5-10(3)

### ■ 分割方法2

プログラムの分かり易さなどの理由で、『D0 I=N1, N2』をそのまま残しておきたい場合などに、図4-5-10(4)の方法を用います(5-3節参照)。まず①で、ループ反復Iとそれを担当するランク値との対応表(MAP)を1回だけ作成しておきます。②ではD0ループを元のプログラムと同じように反復させます。③のIF文で、①で作成した対応表を用いて現在の反復を担当するランクを調べ、それが自分のランク(MYRANK)と一致しているときのみ④の計算を行います。

なお、④の「計算部分」内にさらにD0ループがある場合、通常③のIF文のオーバーヘッドは無視できませんが、④内にD0ループがなく、いきなりステートメントがある場合は、③のIF文によってパフォーマンスが低下してしまうので、この方法は用いない方が良いでしょう。

図4-5-10(4)で配列MAPを確保するのが面倒な場合、図4-5-10(5)のようにもできます。

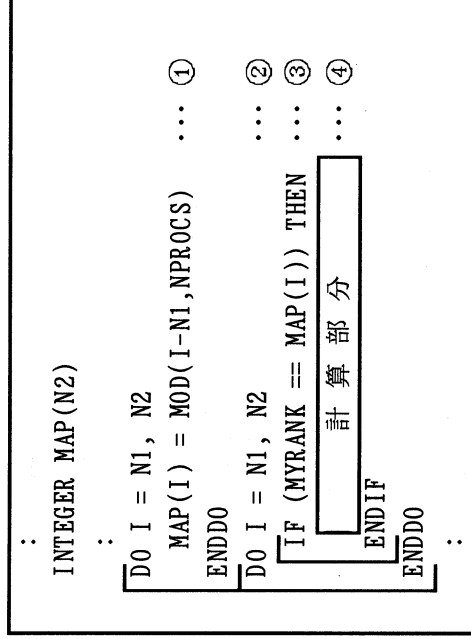


図4-5-10(4)

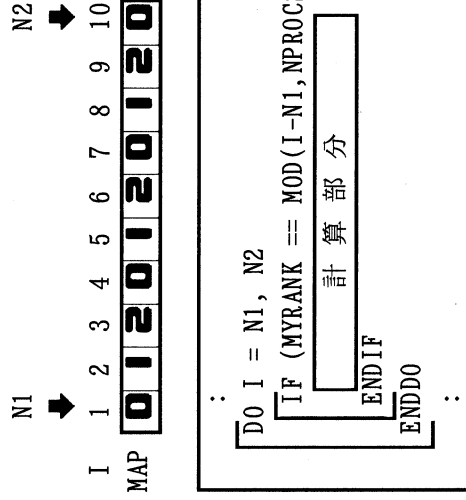


図4-5-10(5)

### ■ 分割方法3

多重ループの内側でサイクリック分割したい場合、以下のように例えば 0→1→2→0→1→2→0→1→2... と反復するカウンター (IRANK) を使用します。

```

:
IRANK = -1
DO J = J1, J2
DO I = I1, I2
  IRANK = IRANK + 1
  IF (IRANK == NPROCS) IRANK = 0
  IF (MYRANK == IRANK) THEN
    計算部分
  ENDIF
ENDDO
ENDDO
:

```

図4-5-10(6)

### ■ 実際の適用例

前節の図4-5-9(1)と同じ例を「分割方法1」でサイクリック分割した例です。

```

PROGRAM MAIN
PARAMETER (N=1000)
DIMENSION A(N)
DO I = 1, N
  A(I) = I
ENDDO
SUM = 0.0
DO I = 1, N
  SUM = SUM + A(I)
ENDDO
PRINT *, 'SUM = ', SUM
END

```

図4-5-11(1)

```

PROGRAM MAIN
INCLUDE 'mpif.h'
PARAMETER (N=1000)
DIMENSION A(N)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
DO I = 1+MYRANK, N, NPROCS ← サイクリック分割
  A(I) = I
ENDDO
SUM = 0.0
DO I = 1+MYRANK, N, NPROCS ← サイクリック分割
  SUM = SUM + A(I)
ENDDO
CALL MPI_REDUCE(SUM, SSUM, 1, MPI_REAL,
& MPI_SUM, 0, MPI_COMM_WORLD, IERR)
SUM = SSUM
IF (MYRANK == 0) PRINT *, 'SUM = ', SUM
CALL MPI_FINALIZE(IERR)
END

```

図4-5-11(2)

# 4-5-6 ブロック・サイクルリック分割

4-5-2節で説明したように、「本当はサイクルリック分割にしたいが、サイクルリック分割にすると別の問題が発生するため、多少サイクルリック性を犠牲にする代わりにその問題を改善する」場合に使用します。本節では、ブロック・サイクルリック分割の具体的な方法を説明します。

## ■ 分割方法1

図4-5-12(1)のD0ループを、プロセス数NPROCSでブロック・サイクルリック分割する方法を図4-5-12(2)に、3プロセスでの分割の様子を図4-5-12(3)に示します。この方法は前節の「分割方法1」に対応します。長さIBLOCK(本例では2)のブロックを1つの単位としてサイクルリック分割を行います。図4-5-12(2)の外側のD0 IIループは図4-5-12(3)の $\Rightarrow$ に対応し、内側のD0 Iループは $\hookrightarrow$ に対応します。

```

:
DO I = N1, N2
  [計算部分]
ENDDO
:

```

図4-5-12(1)元のプログラム

```

:
IBLOCK = ~ ← ブロックの大きさを指定します。
DO II = N1+MYRANK*IBLOCK, N2, NPROCS*IBLOCK
  DO I = II, MIN(II+IBLOCK-1, N2)
    [計算部分]
  ENDDO
ENDDO
:

```

図4-5-12(2)

## ■ 分割方法2

この方法は前節の「分割方法2」に対応し、あらかじめ右図のような対応表(MAP)を作成します。図4-5-12(4)で配列MAPを確保したくない場合、図4-5-12(5)のようにすることもできます。なお、前述のように、IF文を挿入することによってパフォーマンスが低下する可能性がありますので注意して下さい。

```

:
INTEGER MAP(N1:N2)
:
IBLOCK = ~ ← ブロックの大きさを指定します。
DO I = N1, N2
  MAP(I) = MOD(I-N1, NPROCS*IBLOCK)/IBLOCK
ENDDO
DO I = N1, N2
  IF (MYRANK == MAP(I)) THEN
    [計算部分]
  ENDIF
ENDDO
:

```

図4-5-12(4)

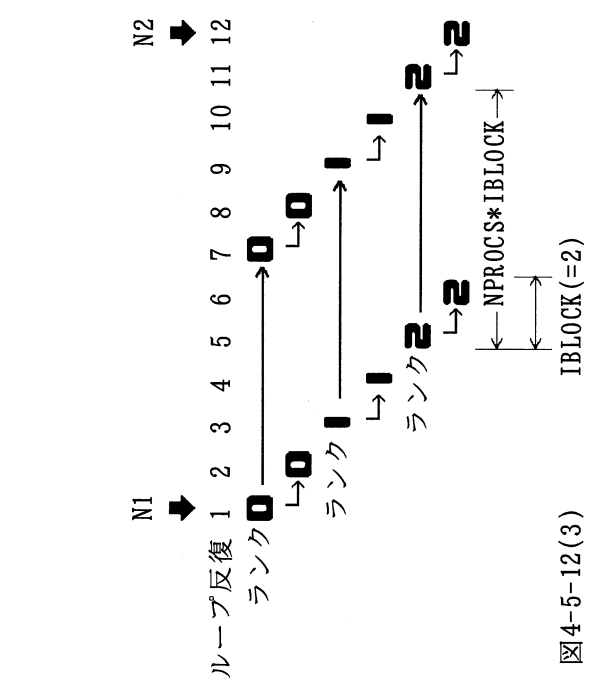


図4-5-12(3)

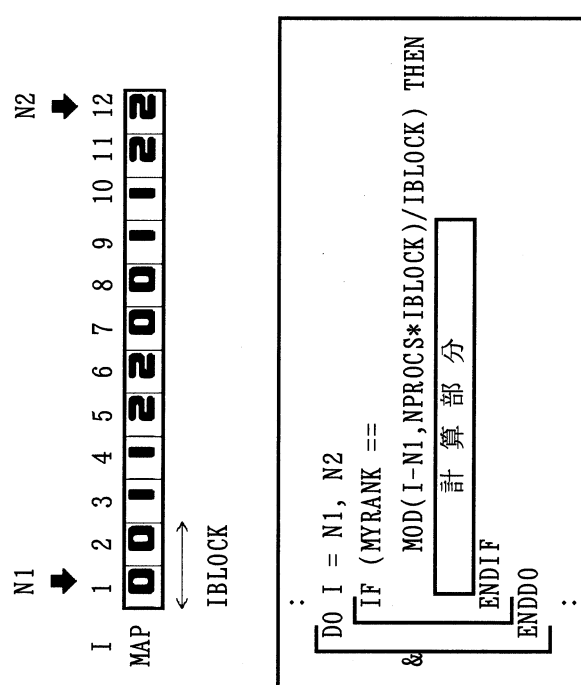


図4-5-12(5)

### ■ 分割方法3

図4-5-12(6)の最初のループで、図4-5-12(7)に示すように、配列MAPにブロックサイクリック分割で自プロセスが担当するループ反復の番号を入れます。最初のループが終了すると、自プロセスが担当する全ループ反復数がICOUNTに入ります。

図4-5-12(6)の2つ目のループで、MAPとICOUNTを使用して、ブロックサイクリック分割で計算を行います。

```

:
INTEGER MAP(*)
:
IBLOCK = ~ ◀ ブロックの大きさを指定します。
ICOUNT = 0
DO I = N1, N2
  IF (MYRANK ==
    MOD(I-N1,NPROCS*IBLOCK)/IBLOCK) THEN
    ICOUNT = ICOUNT + 1
    MAP(ICOUNT) = I
  ENDIF
ENDDO
DO II = 1, ICOUNT
  I = MAP(II)
  ◻ 計算部分
ENDDO
:

```

図4-5-12(6)

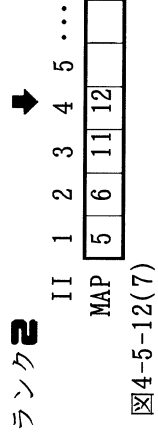
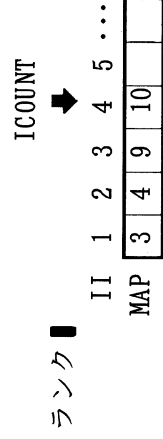
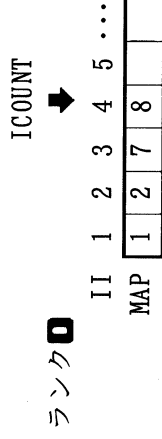
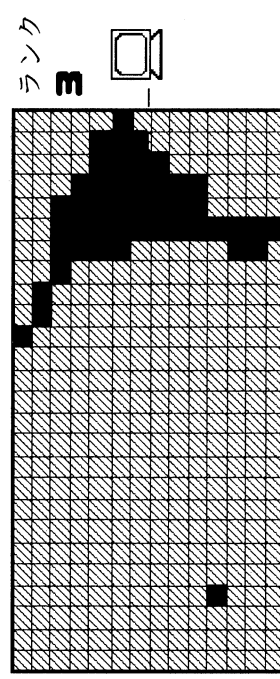
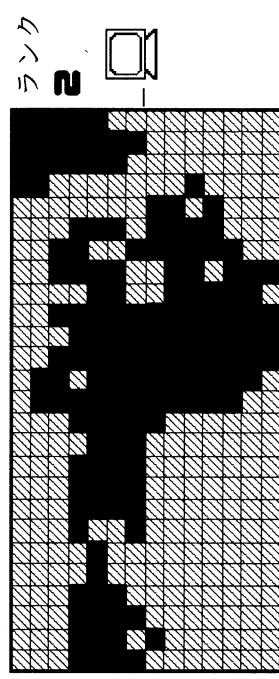
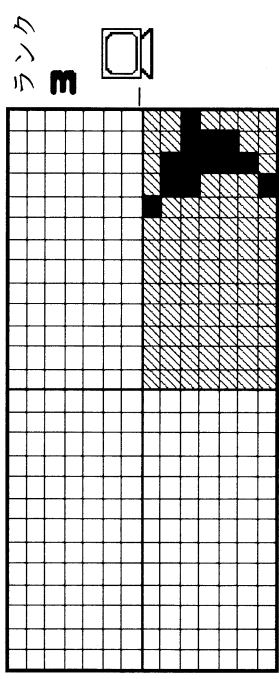
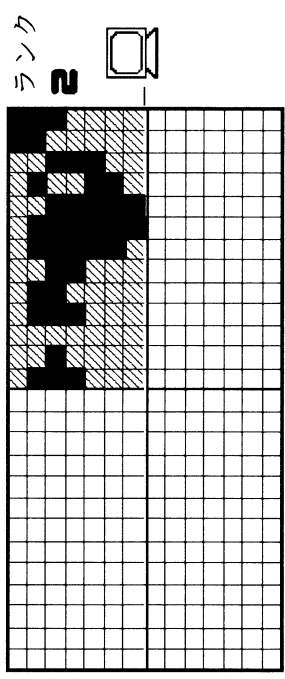
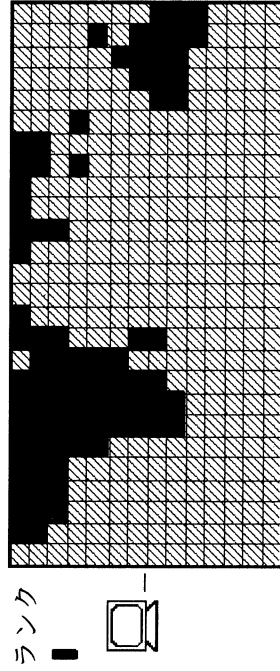
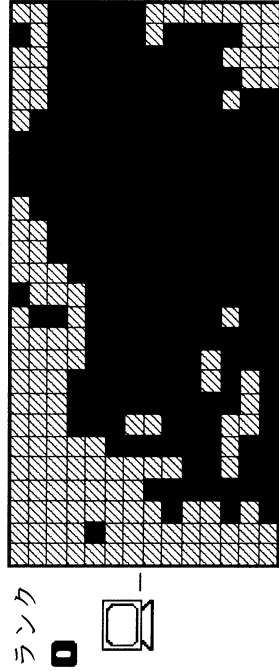
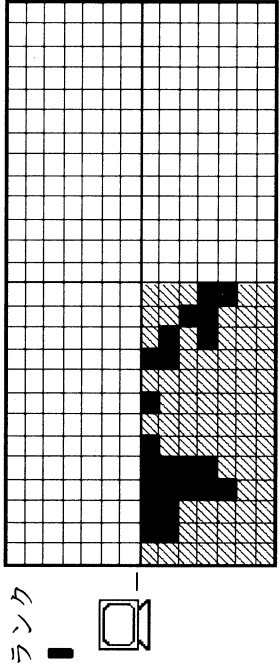
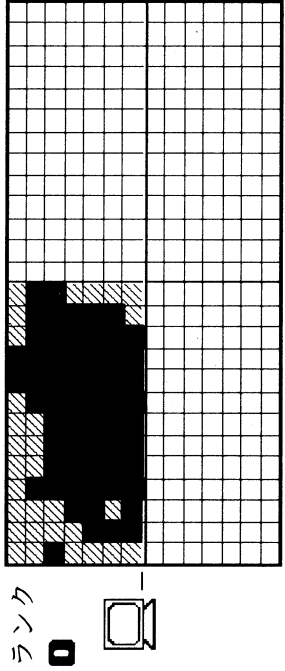
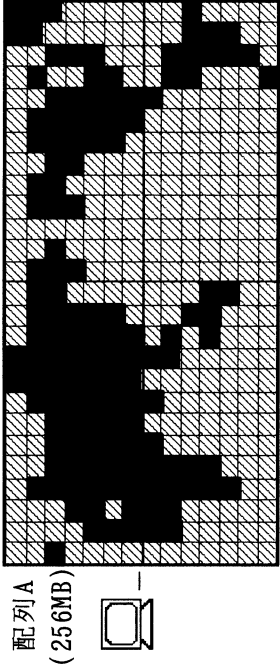


図4-5-12(7)

図4-5-13(1)に示す格子(地球)の大気の流れを計算するプログラムを単体ジョブとして実行するとします。ノード(図)あたり256MBのメモリが搭載されており、図4-5-13(1)は256MBの配列Aを表します。このプログラムを今まで説明した方法で並列化し、4ノードで実行した場合、図4-5-13(2)に示すように、各プロセスが使用するのは配列Aの1/4の部分(図中の着色した部分)のみで、それ以外の部分は使用されません。この並列プログラムを、図4-5-13(3)のように、各プロセスが配列A全体を利用するように改造すれば、全格子数は図4-5-13(1)の4倍となり、よりきめの細かい解析ができるようになります。このように、配列を無駄なく使用することを本書では配列の縮小と呼びます。本節では配列を縮小する方法について説明します。



## 4-5-7-1 配列の縮小方法

### ■ 分散メモリー型並列計算機とメモリー

分散メモリー型並列計算機では、ノード数を増やせばいくらでもメモリーの合計を増やすことができます。従って、配列を縮小すれば、1つの並列ジョブで非常に大きなメモリーを使用することができます。

ところで、配列の縮小は自動的に行われるのではなく、ユーザーがそのようにプログラムを作成すれば、配列の縮小を行うことができます。プログラムによっては配列の縮小は非常に面倒だったり、ほとんど不可能な場合もあります。

一般に、配列を縮小することを意識して最初から並列プログラムを作成する場合はそれほど問題ありませんが、既存の単体プログラムを並列化して配列を縮小する場合は、まず配列を縮小せずに並列化し、答えが合うようになつてから配列の縮小を行うのがよいでしょう。

なお、配列の縮小とは関係ありませんが、プログラム内で使用する大きな配列のうち、メモリー上に常に存在する必要がある配列(作業配列)については、Fortran90の動的割り当て(参考文献[6]の2-6節参照)の機能を使用し、その配列が必要になつた時点でメモリー上に割り当て(ALLOCATE文)、不要になつたらメモリー上から削除する(DEALLOCATE文)ようにすれば、メモリーを節約することができます。

### ■ ブロック分割による配列の縮小方法1

それでは配列の縮小の方法について説明します。例によって図4-5-15(1)の配列の合計を求めるプログラムで説明します。一般化のため、D0ループの範囲がN1からN2までになっていることに注意して下さい。

これを、まず4-5-4節のブロック分割で並列化したのが図4-5-15(2)です。この時点ではまだ配列は縮小されず、例えばプロセス数が3つの場合、各プロセスは配列Aのうち図4-5-14の『配列縮小前』部分を担当しています。

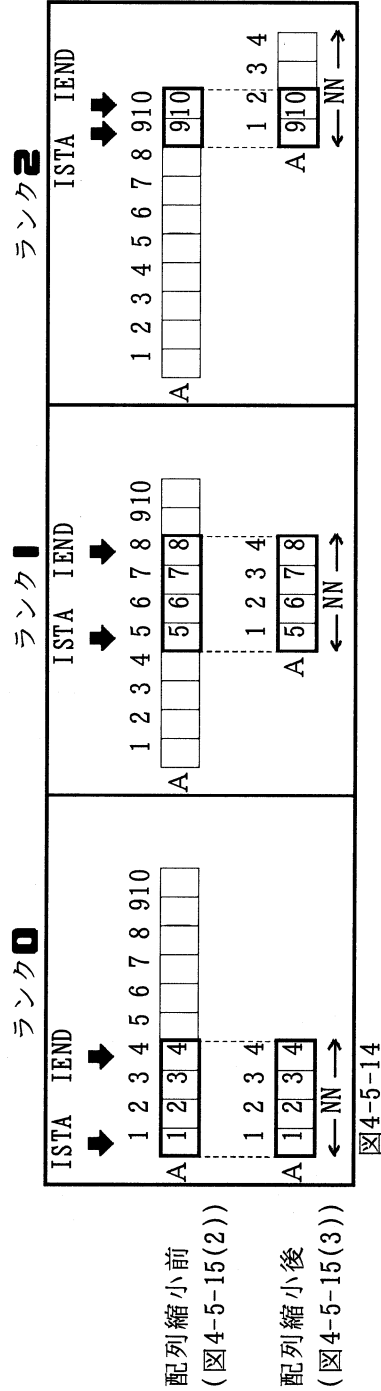


図4-5-14

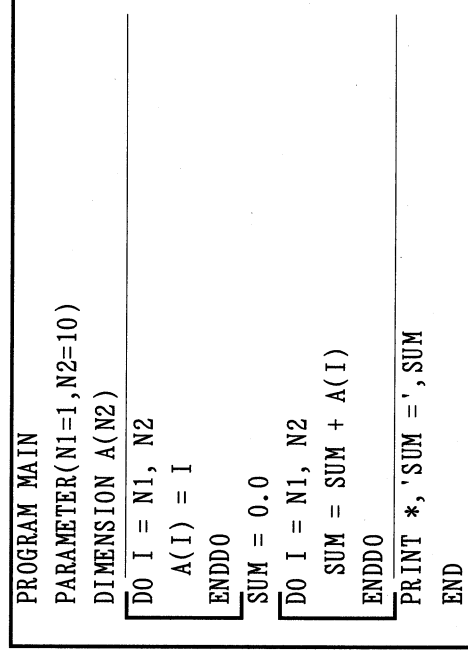


図4-5-15(1)



図4-5-15(2)に対して配列の縮小を行なうと図4-5-15(3)になります。図の『◀』は図4-5-15(2)と異なる部分を示します。図4-5-15(3)について以下に説明します。

● 例えば大きさ100の配列を縮小した場合、ノード数が2ならば 100/2=50、ノード数が4ならば 100/4=25 のように、ノード数によって縮小する配列の大きさを変えする必要があります。従って、縮小する配列Aを宣言する際にノード数がわかって縮小している必要があるため、⑤でノード数(NCPU)を指定します。⑥と⑦で配列Aを縮小して宣言しますが、このときNCPUの値を使用します。なお、⑥では、縮小する前の配列の大きさがノード数で割りきれない場合を考慮しています。

● ⑧では念のために、⑤で指定したノード数(NCPU)と『mpirun -np 3 a.out ▶』で指定したプロセス数(NPROCS)が一致しているかどうかのチェックを行います。

● ③と⑭では、図4-5-15(2)の①、③と同じ反復(ISTAとIEND)をそのまま使用することになります。

● ⑩の左辺と⑫の右辺の配列Aの添字では、ループ反復Iが図4-5-14の『配列縮小後』になるようにアドレス変換を行っています。例えば⑨と⑭のIがISTAの場合、⑩の左辺と⑫の右辺の配列Aの添字は『ISTA-ISTA+1』、つまり『1』となります。

● 図4-5-15(2)の②の右辺では、ループ反復①のIという変数を直接使用しています。一方⑨で①と同じ反復の範囲を使用しているため、⑩の右辺でもIをそのまま使用します。

<pre> PROGRAM MAIN INCLUDE 'mpif.h' PARAMETER(N1=1,N2=10) DIMENSION A(N2) CALL MPI_INIT(IERR) CALL MPI_COMM_SIZE(~_WORLD,NPROCS,IERR) CALL MPI_COMM_RANK(~_WORLD,MYRANK,IERR) CALL PARA_RANGE &amp; (N1,N2,NPROCS,MYRANK,ISTA,IEND) DO I = ISTA, IEND   A(I) = I ENDDO SUM = 0.0 DO I = ISTA, IEND   SUM = SUM + A(I) ENDDO CALL MPI_REDUCE(SUM,SSUM,1,MPI_REAL, &amp; MPI_SUM,0,MPI_COMM_WORLD,IERR) SUM = SSUM IF (MYRANK==0) PRINT *, 'SUM = ',SUM CALL MPI_FINALIZE(IERR) END </pre>	<pre> PROGRAM MAIN INCLUDE 'mpif.h' PARAMETER(N1=1,N2=10) PARAMETER(NCPU=3) PARAMETER(NN=(N2-N1)/NCPU+1) DIMENSION A(NN) CALL MPI_INIT(IERR) CALL MPI_COMM_SIZE(~_WORLD,NPROCS,IERR) CALL MPI_COMM_RANK(~_WORLD,MYRANK,IERR) IF (NPROCS /= NCPU) THEN   IF (MYRANK==0) PRINT *, '===ERROR==='   CALL MPI_FINALIZE(IERR) STOP ENDIF CALL PARA_RANGE &amp; (N1,N2,NPROCS,MYRANK,ISTA,IEND) DO I = ISTA, IEND   A(I-ISTA+1) = I ENDDO SUM = 0.0 DO I = ISTA, IEND   SUM = SUM + A(I-ISTA+1) ENDDO CALL MPI_REDUCE(SUM,SSUM,1,MPI_REAL, &amp; MPI_SUM,0,MPI_COMM_WORLD,IERR) SUM = SSUM IF (MYRANK==0) PRINT *, 'SUM = ',SUM CALL MPI_FINALIZE(IERR) END </pre>
--	---

図4-5-15(3)

図4-5-15(2)

## ■ ブロック分割による配列の縮小方法2

図4-5-15(3)では、⑨と⑩のループ反復は図4-5-15(2)の①、③をそのまま使用しました。ところで、配列Aを縮小した場合、図4-5-14の『配列縮小後』に示すように、各プロセスは配列Aを『1』から使用します。これにともない、⑨と⑩のループ反復を『1』から開始させるという方法もあります。これを図4-5-15(4)に示します。図4-5-15(4)の『◀』は図4-5-15(3)と異なる部分を示します。

● ⑬と⑭に示すように、ループ反復の下限を『1』から開始しています。ループ反復の上限は、⑨と⑩からわかるように、『IEND-ISTA+1』となります。ループ反復の下限と上限を変更したため、⑭の左辺と⑯の右辺の配列Aの添字はIとなります。これは元の②、④と同じです。

● ⑩では②の右辺のIをそのまま使用しましたが、今度は⑬の反復を『1』から開始するようにしたため、Iを⑭の右辺のように変更する必要があります。⑬でI=1のとき⑭の右辺は『1+ISTA-1』、つまりISTAとなりません。⑨でI=ISTAのとき⑩の右辺もISTAとなるので、両者は同じであることがわかります。

図4-5-15(3)と(4)のどちらの形式を使用するかは好みによります。プログラムをまず縮小せずに並列化し、その後で縮小する場合は、図4-5-15(3)の方がD0ループの反復を修正する必要があるため簡単でしょう。逆に、配列を縮小した並列化を最初から意識してプログラムを作成する場合は、図4-5-15(4)の方が自然かもしれません。なお、⑩と⑫の『-STA+1』、⑬と⑮の『IEND-STA+1』などは定数にするともう少しすっきりします。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER(N1=1,N2=10)
  PARAMETER(NCPU=3)
  PARAMETER(NN=(N2-N1)/NCPU+1)
  DIMENSION A(NN)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~_WORLD,NPROCS,IERR)
  CALL MPI_COMM_RANK(~_WORLD,MYRANK,IERR)
  IF (NPROCS /= NCPU) THEN
    IF (MYRANK==0) PRINT *, '===ERROR=== '
    CALL MPI_FINALIZE(IERR)
    STOP
  ENDF
  CALL PARA_RANGE
  & (N1,N2,NPROCS,MYRANK,ISTA,IEND)
  DO I = 1, IEND-ISTA+1
    A(I) = I+ISTA-1
  ENDDO
  SUM = 0.0
  DO I = 1, IEND-ISTA+1
    SUM = SUM + A(I)
  ENDDO
  CALL MPI_REDUCE(SUM,SSUM,1,MPI_REAL,
  & MPI_SUM,0,MPI_COMM_WORLD,IERR)
  SUM = SSUM
  IF (MYRANK==0) PRINT *, 'SUM = ', SUM
  CALL MPI_FINALIZE(IERR)
  END

```

図4-5-15(4)

■ ブロック分割による配列の縮小方法3

Fortran90の新機能であるメモリの動的割振り機能(C言語の関数mallocに相当)を用いると、配列の縮小を比較的簡単に行うことができます。まず図4-5-16(1)の例で、メモリの動的割振りの方法について簡単に説明します。なお、メモリの動的割振りの詳細は、参考文献[6]の2-6節を参照して下さい。

- ①の下線部で、REAL\*8の配列AとBがプログラムの実行時に動的に割振られることを宣言します。配列が1次元の場合はA(:)、2次元の場合はB(:,:)のように指定します。下線部の途中のブランクはなくてもかまいません。
- この例では、②で標準入力から配列AとBの大きさIMAX, JMAXを読み込みます。読み込んだIMAX, JMAXを使用して③を実行すると、配列AとBが実際にメモリー上に割振られます。以後は配列AとBを通常の配列と同様に使用することができます。なお③の下線部(任意)を指定すると、整数IERR(名前は任意)に、ALLOCATE文が成功したときはゼロが、(メモリー不足などで)失敗したときはゼロ以外の値が戻ります。これを④でチェックし、エラーの場合は終了します。
- 配列AとBを引数でサブルーチンに渡す場合は④のようになります。
- 使用が終了して不要になった配列を解放したい場合は⑤を実行します。もちろん解放せずにプログラムを終了してもかまいません。⑤の下線部は任意です。
- 配列AとBを図4-5-16(2)の⑥のようにしてサブルーチンに渡すこともできます。

```

PROGRAM MAIN
REAL*8, ALLOCATABLE :: A(:),B(:,:) ①
READ(5,*) IMAX, JMAX ②
ALLOCATE(A(IMAX),B(IMAX, JMAX), STAT=IERR) ③
IERRが0以外だったらエラーで終了する。 ③
CALL SUB(A, B, IMAX, JMAX) ④
DEALLOCATE(A, B, STAT=IERR) ⑤
IERRが0以外だったらエラーで終了する。
END
SUBROUTINE SUB(A, B, IMAX, JMAX) ④
REAL*8 A(IMAX), B(IMAX, JMAX) ④
:
END
    
```

図4-5-16(1)

```

PROGRAM MAIN
COMMON /COM/ IMAX, JMAX ⑥
REAL*8, ALLOCATABLE :: A(:),B(:,:)
READ(5,*) IMAX, JMAX
ALLOCATE(A(IMAX), B(IMAX, JMAX), STAT=IERR)
IERRが0以外だったらエラーで終了する。 ⑥
CALL SUB(A, B)
DEALLOCATE(A, B, STAT=IERR)
IERRが0以外だったらエラーで終了する。
END
SUBROUTINE SUB(A, B) ⑥
COMMON /COM/ IMAX, JMAX ⑥
REAL*8 A(IMAX), B(IMAX, JMAX)
:
END
    
```

図4-5-16(2)

次に、動的に割振られた配列を引数でなくCOMMON文でサブルーチンに渡す方法を説明します。実は動的に割振られた配列をCOMMON文の中に指定することはできません。その代わりに、Fortran90の新機能であるMODULE文とUSE文を使用します。MODULE文はCOMMON文を含んだINCLUDEファイルに相当し、USE文はINCLUDE文に相当します。図4-5-16(3)の⑦に示すようにMODULE文に適切な名前(本例ではCOM)をつけます。MODULE～ENDで囲まれた⑧の部分には、グローバルとして扱いたい変数や配列を指定します。この部分に、動的割振りする変数や配列を指定することができます。⑧で指定した配列や変数は、⑨のUSE文を指定したメイン(サブ)ルーチンのみから参照/更新することができま

す。MODULE文とUSE文に関する注意点を以下に述べますが、使用法に少くせがあるので注意して下さい(詳細は参考文献[6]の2-6節参照)。

- USE文はメイン(サブ)ルーチンの先頭で指定する必要があります。
- MODULE～ENDはINCLUDEファイルと異なりプログラムと同じファイルに入れることができますが、ファイル内の最初のUSE文よりも前に置かれなければなりません。従ってファイル内の先頭に置くのがよいでしょう。
- MODULE～ENDの中にIMPLICIT文を指定することもできますが、これはMODULE～END内のみで有効となり、USE文を指定したメイン(サブ)ルーチン側には適用されないので、必要ならばメイン(サブ)ルーチンで再度IMPLICIT文を宣言して下さい。

なお、この例では⑧でIMAX, JMAXを指定していますが、これは動的割振りとは関係なく、サブルーチンSUBでIMAXとJMAXを使用しなければこの指定は不要です。

MODULE文はCOMMON文を含んだINCLUDEファイルに相当しますが、INCLUDEファイルと異なりプログラムの入ったファイルと別ファイルにする必要はありません。

```

MODULE COM
⑦
REAL*8, ALLOCATABLE :: A(:, :), B(:, :)
⑧
INTEGER IMAX, JMAX
⑧
END
PROGRAM MAIN
⑨
USE COM
:
READ(5, *) IMAX, JMAX
ALLOCATE(A(IMAX), B(IMAX, JMAX), STAT=IERR)
IERRが0以外だったらエラーで終了する。
:
CALL SUB
:
DEALLOCATE(A, B, STAT=IERR)
END
SUBROUTINE SUB
⑨
USE COM
:
END

```

図4-5-16(3)

```

PROGRAM MAIN
INCLUDE 'mpif.h'
PARAMETER(N1=1, N2=10)
REAL, ALLOCATABLE :: A(:)
①
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(~_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK(~_WORLD, MYRANK, IERR)
CALL PARA_RANGE
& (N1, N2, NPROCS, MYRANK, ISTA, IEND)
②
ALLOCATE(A(ISTA:IEND), STAT=IERR)
③
IERRが0以外だったらエラーで終了する。
DO I = ISTA, IEND
A(I) = I
ENDDO
SUM = 0.0
DO I = ISTA, IEND
SUM = SUM + A(I)
ENDDO
CALL MPI_REDUCE(SUM, SSUM, 1, MPI_REAL,
& MPI_SUM, 0, MPI_COMM_WORLD, IERR)
SUM = SSUM
IF (MYRANK==0) PRINT *, 'SUM = ', SUM
CALL MPI_FINALIZE(IERR)
END

```

図4-5-17(1)

メモリーの動的割振り機能を用いて図4-5-15(1)の配列を縮小した例を上記の図4-5-17(1)に示します。まず①で配列Aが動的に割振られることを宣言します。②で自分の担当範囲の下限(ISTA)と上限(IEND)を確定した後、③で自分の担当部分の1次元配列A(ISTA:IEND)をメモリー上に割振ります。これで全て完了です。この方法では図4-5-14は図4-5-17(2)のようになります。プロセスによって配列Aの下限と上限、大きさが異なっていることに注意して下さい。

この方法では、図4-5-15(3)の⑤⑥⑧は全て不要になり、また⑩や⑭のように配列の添字を修正したり、図4-5-15(4)の⑮や⑯のようにD0ループの反復を修正する必要もなくなります。

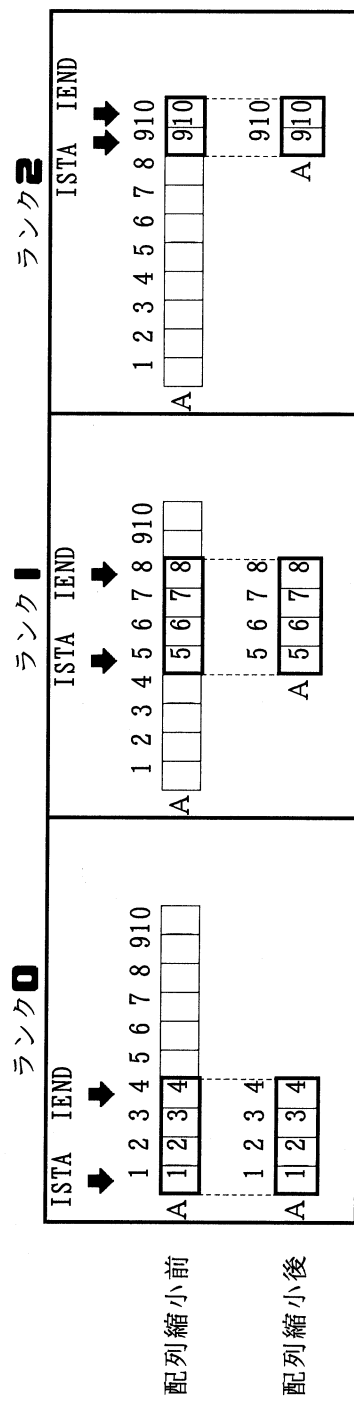


図4-5-17(2)

■ 配列を縮小するメリット

通常、配列を縮小するのは、1台のノードが持っているメモリーだけではメモリーが足りない場合があります。逆に言うと、1台のノードが持っているメモリーで足りるのなら、手間をかけて配列を縮小する必要はないわけですが、ところがこのような場合でも、計算と通信のパフォーマンスの点から、やはり配列を縮小した方がよい場合があります。

図4-5-18(1)で、縦方向が配列の1次元目、すなわちメモリーが連続する方向とします。図4-5-18(1)上段のように、配列を1次元目でブロック分割し縮小をしない場合、ランク0のプロセスが処理する斜線の部分のうち、①と②の要素は↓に示すようにメモリー上で離れています。このため、ランク0のプロセスが斜線の部分の要素を順に計算する場合、キャッシュ・ミスが発生してパフォーマンスが若干低下します。また、斜線の部分を通信する際にもこの部分で同様にキャッシュ・ミスが発生してパフォーマンスが若干低下します。

これを図4-5-18(1)下段のように配列を縮小すれば、メモリー上で①と②が連続するため、これらのキャッシュ・ミスがなくなります。なお、配列の2次元目でブロック分割した場合は、縮小しなくても全ての要素がメモリー上で連続しているため、この問題はありませぬ。

図4-5-18(2)上段のように、配列をサイクリック分割し縮小をしない場合も、メモリー上で③と④が離れているため、図4-5-18(1)と同様の問題が発生します。

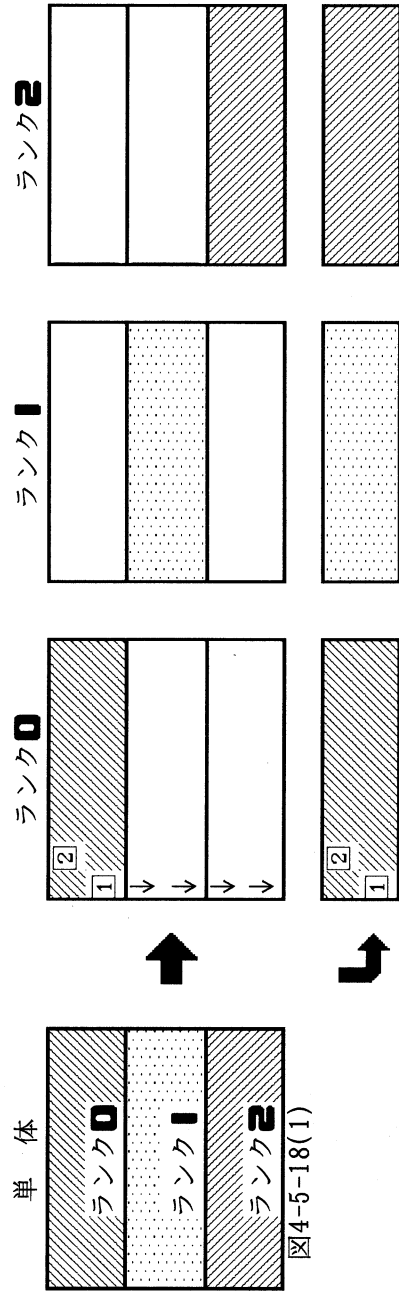


図4-5-18(1)

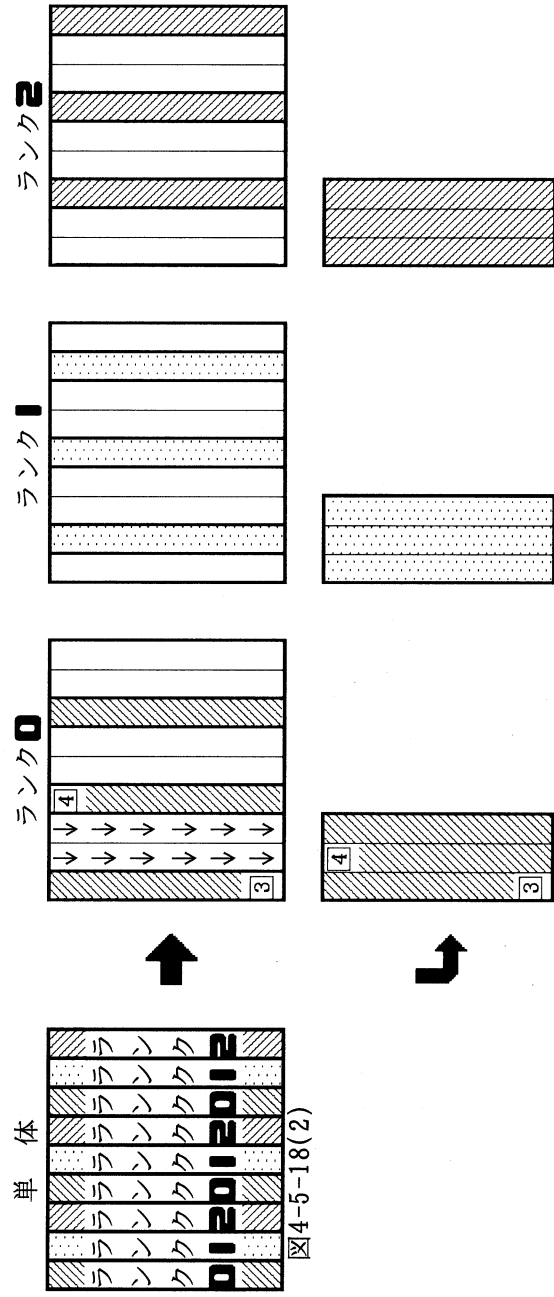


図4-5-18(2)

## 4-5-7-2 配列の縮小に伴う変更(I/O以外の部分)

単体プログラムをまず配列を縮小せずに並列化し、その後で配列を縮小した場合、以下の箇所を修正する必要があります。この他に次節で説明するI/O部分の修正が必要になります(プログラムによってはそれ以外にも修正の必要な箇所があるかも知れません)。

なお、以下の例では4-5-7-1節の「ブロック分割による配列の縮小方法3」で縮小しています。

### ■ サブルーチン化した通信ルーチンの配列宣言の変更

配列を縮小せずに並列化したプログラムを図4-5-20(1)に示します。このプログラムでは異なる配列AとBに対して①で同じパターンでの通信を行うため、②のように通信部分をサブルーチン化しています。このプログラムの配列AとBを図4-5-20(2)の③のように縮小した場合、④の部分も修正する必要がありますがあります。

```

:
REAL A(N),B(N)
ISTA = ~
IEND = ~
CALL SHIFT(A,N,ISTA,IEND) ①
CALL SHIFT(B,N,ISTA,IEND) ①
:
SUBROUTINE SHIFT(X,N,ISTA,IEND) ②
REAL X(N)
CALL MPI_XXXX(X(ISTA),~)
:

```

図4-5-20(1) 並列化(縮小前)

```

:
REAL,ALLOCATABLE::A(:),B(:)
ISTA = ~
IEND = ~
ALLOCATE(A(ISTA:IEND),B(ISTA:IEND)) ③
CALL SHIFT(A,N,ISTA,IEND)
CALL SHIFT(B,N,ISTA,IEND)
:
SUBROUTINE SHIFT(X,N,ISTA,IEND)
REAL X(ISTA:IEND) ④
CALL MPI_XXXX(X(ISTA),~)
:

```

図4-5-20(2) 並列化(縮小後)

### ■ それまで並列化していなかったDOループの並列化

図4-5-21(1)の単体プログラムでは、②のタイムステップ内に含まれる③の部分で計算時間がかかっており、①の部分の計算時間はわずかである⑤の部分のみを並列化することもできません。

一方配列Aを縮小した図4-5-21(3)では、計算時間のかからない⑥の部分も並列化する必要があります。

```

:
DO I=1,N ①
A(I) = 1.0
ENDDO
DO ITIME=1,1000 ②
DO I=1,N ③
A(I) = A(I) + 1.0
ENDDO
:
ENDDO
:

```

図4-5-21(1) 並列化前

```

:
DO I=1,N ④
A(I) = 1.0
ENDDO
DO ITIME=1,1000
DO I=ISTA,IEND ⑤
A(I) = A(I) + 1.0
ENDDO
:
ENDDO
:

```

図4-5-21(2) 並列化(縮小前)

```

:
ALLOCATE(A(ISTA:IEND)) ⑥
DO I=ISTA,IEND
A(I) = 1.0
ENDDO
DO ITIME=1,1000
DO I=ISTA,IEND
A(I) = A(I) + 1.0
ENDDO
:
ENDDO
:

```

図4-5-21(3) 並列化(縮小後)

■ 派生データ型を使用した通信の修正

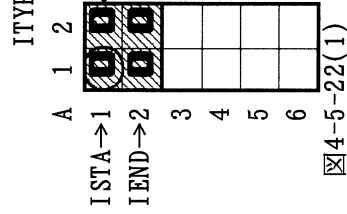
メモリー上で飛び飛びになっているデータを送(受)信する場合、派生データ型(3-5節参照)を使用することができず、そして飛び飛びのデータを含む配列を縮小した場合、データは(一般的に)連続するので、派生データ型を使用した通信部分を修正する必要があります。

図4-5-22(1)で、各プロセスは2次元配列A(6,2)のうち着色した部分のデータを送信するとします。着色した部分に含まれる要素はメモリー上で飛び飛びなので、派生データ型を使用します。図4-5-23(1)の①で、3-5-3-2節で説明した自作のサブローチンPARA\_TYPE\_BLOCK2Aをコールし、図4-5-22(1)の○の要素から始まって着色した部分を表す派生データ型ITYPEを作成し、これを使用して②で通信を行います。

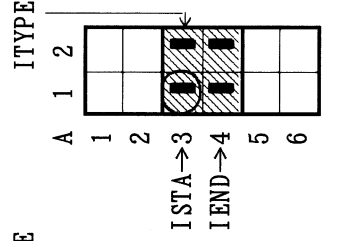
図4-5-23(2)では、③で自作のサブローチンPARA\_TYPE\_BLOCK2コールし(MPI-2が使用できるマシン環境の場合は、MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照)を使用して下さい)、図4-5-22(2)の○から始まって着色した部分を表す派生データ型ITYPEを作成し、これを使用して④で通信を行います。なお④の送信データの先頭アドレスは②と異なります。

この配列Aを図4-5-22(3)のように縮小した場合、送信データはメモリー上で連続するので、通信で派生データ型を使用する必要はありません。図4-5-23(3)の⑤と⑥で縮小した配列Aをアロケートし、⑦の下線部で送信するデータの個数を指定し、(派生データ型を使用せずに)通常の方法でデータを送信します。

【ランク0】



【ランク1】



【ランク2】

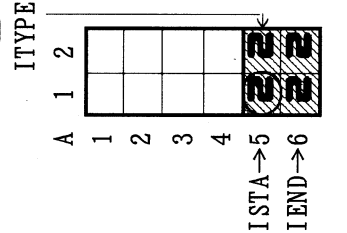


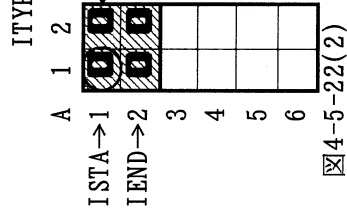
図4-5-22(1)

```

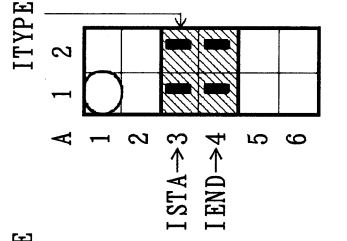
:
REAL A(6,2)
  ISTAとIENDを決定する。
CALL PARA_TYPE_BLOCK2A
& (1,6,IEND-ISTA+1,2,MPI_REAL,ITYPE)
CALL MPI_SEND(A(ISTA,1),1,ITYPE,~)
:
    
```

図4-5-23(1)

【ランク0】



【ランク1】



【ランク2】

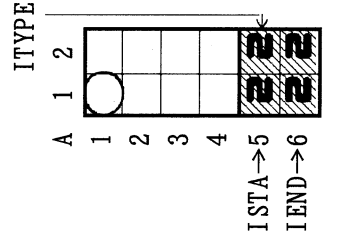


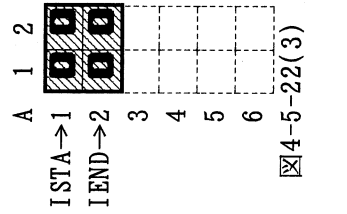
図4-5-22(2)

```

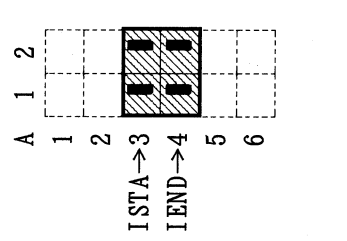
:
REAL A(6,2)
  ISTAとIENDを決定する。
CALL PARA_TYPE_BLOCK2
& (1,6,1,ISTA,IEND,1,2,MPI_REAL,ITYPE)
CALL MPI_SEND(A,1,ITYPE,~)
:
    
```

図4-5-23(2)

【ランク0】



【ランク1】



【ランク2】

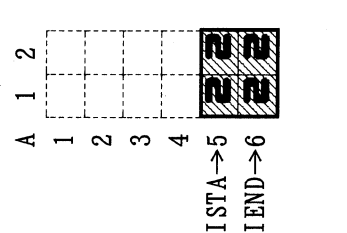


図4-5-22(3)

```

:
REAL,ALLOCATABLE::A(:,:)
  ISTAとIENDを決定する。
ALLOCATE(A(ISTA:IEND,2))
CALL MPI_SEND(A(ISTA,1),(IEND-ISTA+1)*2,
& MPI_REAL,~)
:
    
```

図4-5-23(3)

## 4-5-7-3 配列の縮小に伴う変更(I/O部分)

本節では、並列プログラム内でI/Oが行われる配列を縮小した場合の修正方法について説明します。

図4-5-24(1)の単体プログラムでは、①で配列X(2,6)を計算し、②で1回のWRITE命令で10番ファイルに書き出します。また図4-5-24(2)の単体プログラムでは、③で10番ファイルを1回のREAD命令で配列Y(2,6)に読み込んだ後、④で計算を行います。この2つのプログラムを並列化し、配列XとYを縮小するとします。

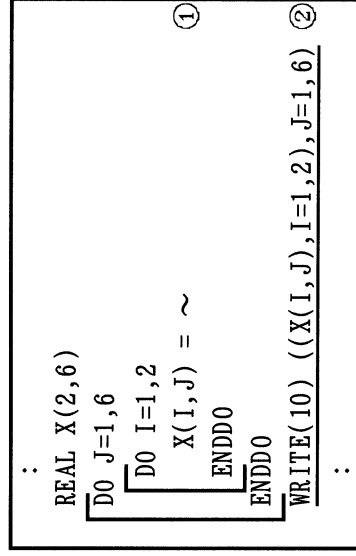


図4-5-24(1) 出力プログラム

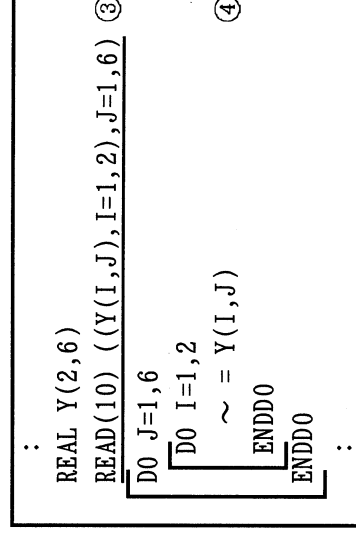
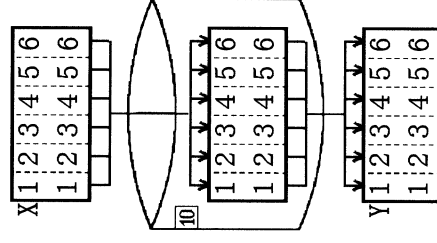


図4-5-24(2) 入力プログラム



まず図4-5-24(1)の並列化ですが、4-4-2節の出力パターンのうち、「出力パターン1」を採用し、計算が終了した後で合体しなければならぬので面倒です。そこで「出力パターン1」を採用します。なお、「出力パターン3」のMPI-I0の機能を使用することも可能ですが、説明は省略します。

配列を縮小する場合は(一般に)メモリーに余裕がないので、各プロセスから送信されたデータ全体を同時にランク0の受信用配列に集めることはできません。そこでランク0のプロセスは、まずランク1のプロセスから送信されたデータを受信するために必要な最小限の大きさの一時配列tempを確保してから受信を行い、それを出力ファイルに書き出し、不要になった配列tempを解放します。そして次にランク2に対して同じ処理を行います(プロセス数が3の場合)。

図4-5-24(1)を並列化し、配列Xを縮小したプログラムとその動作を図4-5-25(1)(2)に示します。なお、図4-5-25(1)の出力ファイルは図4-5-24(2)のプログラムで読むことはできません(レコード形式が異なるので)。読む場合は図4-5-25(3)(単体プログラム)または図4-5-28(1)(並列プログラム)を使用して下さい。

- 図4-5-25(1)の[1]で、配列X(2,6)の2次元目を、4-5-4節で説明したサブルーチンPARA\_RANGEを使用してブロック分割し、各プロセスが担当する2次元目の下限と上限を配列JJSTA, JJENDに設定します。
- [2]で各プロセスは、自分が担当する2次元目の下限と上限をJSTA, JENDに設定します。
- [3]で各プロセスは、自分が担当する大きさの配列Xを動的割振りで縮小して確保し、[4]で計算を行います。
- [5]でランク0のプロセスは、まず自分が担当した部分をファイルに書き出します(図4-5-25(2)の[5]参照)。図4-5-24(1)では配列X全体を1回のWRITE文で書き出しましたが、[5]では1列ずつ書き出します。この理由は(注)を参照して下さい。
- [6]のDOループでまずIRANK=1のとき、ランク0のプロセスはランク1のプロセスから送信されたデータを以下のように書き出します。
  - ランク1(正確には0以外)のプロセスは、[4]で自分の担当部分をランク0のプロセスに送信します。
  - ランク0のプロセスは、[7]で、ランク1のプロセスから送信されたデータを受信するために必要な最小限の大きさの一時配列tempを動的割振りで確保します。次に[8]でデータを受信し、[9]で10番ファイルに列ごとに書き出します。最後に[10]で不要となった一時配列tempを解放します。
- [6]のDOループがIRANK=2となり、ランク2から送信されたデータを上記と同様に受信して書き出します。以後同様に処理を行います。

(注) [5]と[9]で、配列全体を1回のWRITE文で書き出した場合は図4-5-26(1)の下線のようになり、図4-5-26(2)のファイル内の実線で囲んだ部分が1回のWRITEで1レコードとして書き出されます(例えば3プロセスで並列化した場合、3レコードが書き出されます)。すると、後でこのファイルを単体ジョブ、またはプロセス数が3以外の並列ジョブで読み込む場合、読み込みが面倒になります、このため図4-5-25(1)では図4-5-25(2)のように1列ずつバラバラに書き出しました。ただし図4-5-26(1)のように配列全体を書き出すよりも、パフォーマンスは多少低下するかも知れません。



```

:
INCLUDE 'mpif.h'
REAL, ALLOCATABLE::TEMP(:, :)
REAL, ALLOCATABLE::X(:, :)
INTEGER ISTATUS(MPI_STATUS_SIZE)
PARAMETER(NCPU=3)
INTEGER JJSTA(0:NCPU-1), JJEND(0:NCPU-1)
:
DO IRANK=0, NPROCS-1
CALL PARA_RANGE(1,6,NPROCS,IRANK,
& JJSTA(IRANK),JJEND(IRANK))
ENDDO
JSTA = JJSTA(MYRANK)
JEND = JJEND(MYRANK)
ALLOCATE(X(2,JSTA:JEND))
DO J=JSTA, JEND
DO I=1,2
X(I,J) = ~
ENDDO
ENDDO
IF (MYRANK==0) THEN
DO J=JSTA, JEND
WRITE(10)(X(I,J), I=1,2)
ENDDO
DO IRANK=1, NPROCS-1
ALLOCATE
(TEMP(2, JJSTA(IRANK):JJEND(IRANK))) [7]
ISIZE=2*(JJEND(IRANK)-JJSTA(IRANK)+1)
CALL MPI_RECV(TEMP, ISIZE, MPI_REAL,
& IRANK, 1, MPI_COMM_WORLD, IREQ, IERR) [8]
CALL MPI_WAIT(IREQ, ISTATUS, IERR) [8]
DO J=JJSTA(IRANK), JJEND(IRANK)
WRITE(10)(TEMP(I,J), I=1,2)
ENDDO
DEALLOCATE(TEMP)
ENDDO
ELSE
CALL MPI_ISEND(X, 2*(JEND-JSTA+1), MPI_REAL,
& 0, 1, MPI_COMM_WORLD, IREQ, IERR) [11]
CALL MPI_WAIT(IREQ, ISTATUS, IERR)
ENDIF
:

```

図4-5-25(1) 図4-5-25(2)のプログラム例

```

:
REAL Y(2,6)
DO J=1,6
READ(10)(Y(I,J), I=1,2) ③
ENDDO
:

```

図4-5-25(3)

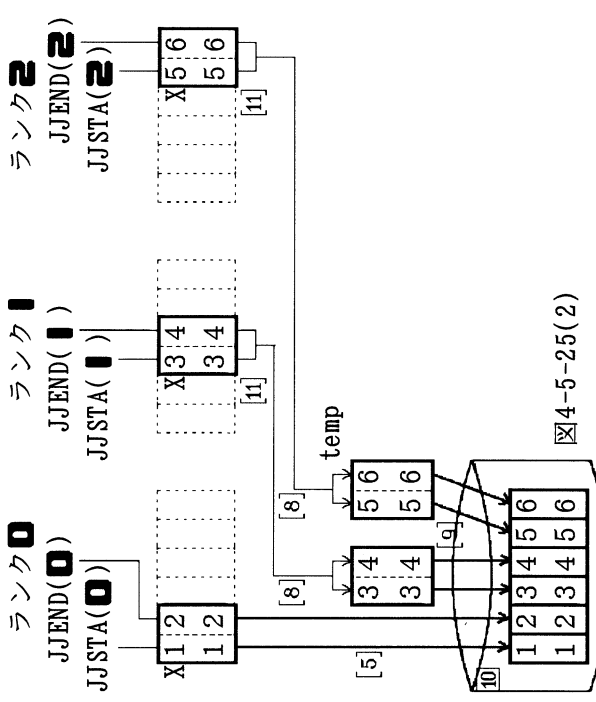


図4-5-25(2)

```

:
INCLUDE 'mpif.h'
REAL, ALLOCATABLE::TEMP(:, :)
:
IF (MYRANK==0) THEN
WRITE(10)((X(I,J), I=1,2), J=JSTA, JEND)
DO IRANK=1, NPROCS-1
ALLOCATE
(TEMP(2, JJSTA(IRANK):JJEND(IRANK)))
ISIZE=2*(JJEND(IRANK)-JJSTA(IRANK)+1)
CALL MPI_RECV(TEMP, ISIZE, MPI_REAL,
& IRANK, 1, MPI_COMM_WORLD, IREQ, IERR)
CALL MPI_WAIT(IREQ, ISTATUS, IERR)
WRITE(10)((TEMP(I,J), I=1,2),
& J=JJSTA(IRANK), JJEND(IRANK))
DEALLOCATE(TEMP)
ENDDO
:

```

図4-5-26(1) 図4-5-26(2)のプログラム例

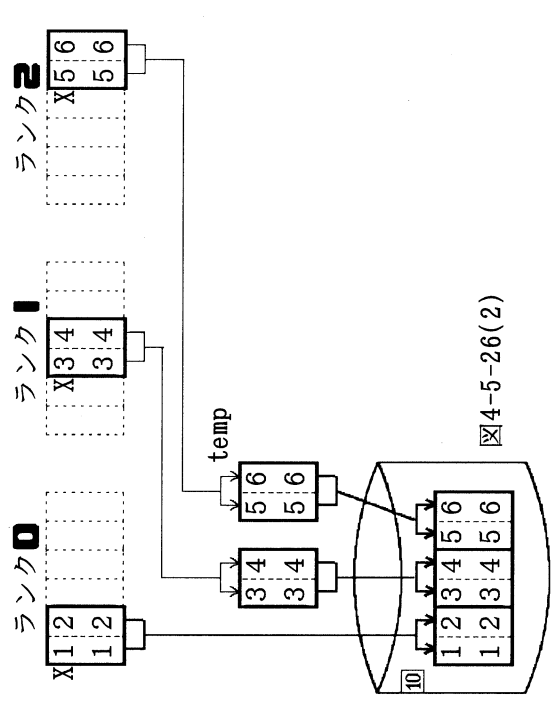


図4-5-26(2)

図4-5-25(1)では、1つのプロセスから送信されたデータが全て入る大きさの一時配列tempで受信しましたが、ファイルへの書き出しは1列ずつ行いました。もし一時配列tempの大きさも節約したいのであれば、図4-5-27(1)(2)に示すように、配列tempを1列分の大きさとし、各プロセスからのデータを1列ずつ受信してはファイルに書き出すこともできます。ただし通信回数が増えるので、パフォーマンスは若干低下します。

最後に、図4-5-24(2)の入力プログラムを4-4-1節の「入力パターン3」で並列化する方法を図4-5-28(1)(2)に示します。入力ファイルは図4-5-25(1)または図4-5-27(1)で作成したファイルであるとしします。

- 図4-5-28(1)の[1]で2次元目をブロック分割し、[2]で自分が担当する2次元目の下限と上限をJSTA, JENDに設定します。[3]で配列Yを縮小して動的割振りで確保します。
- [4]で入力ファイルのうち自分の担当範囲より前の部分を1列ずつ空読みします。
- [5]で入力ファイルのうち自分の担当範囲の部分を配列Yに1列ずつ読み込みます。
- [6]で自分の担当範囲の部分を計算します。

```

:
INCLUDE 'mpif.h'
REAL TEMP(2)
!
! 図4-5-25(1)のAと同じ
!
IF (MYRANK=0) THEN
DO J=JSTA,JEND
WRITE(10)(X(I,J), I=1,2)
ENDDO
DO IRANK=1, NPROCS-1
DO J=JJSTA(IRANK), JJEND(IRANK)
CALL MPI_RECV(TEMP, 2, MPI_REAL,
IRANK, 1, MPI_COMM_WORLD, IREQ, IERR)
CALL MPI_WAIT(IREQ, ISTATUS, IERR)
WRITE(10)(TEMP(I), I=1,2)
ENDDO
ENDDO
ELSE
DO J=JSTA,JEND
CALL MPI_SEND(X(1,J), 2, MPI_REAL,
0, 1, MPI_COMM_WORLD, IREQ, IERR)
CALL MPI_WAIT(IREQ, ISTATUS, IERR)
ENDDO
ENDIF
:
&

```

図4-5-27(1) 図4-5-27(2)のプログラム例

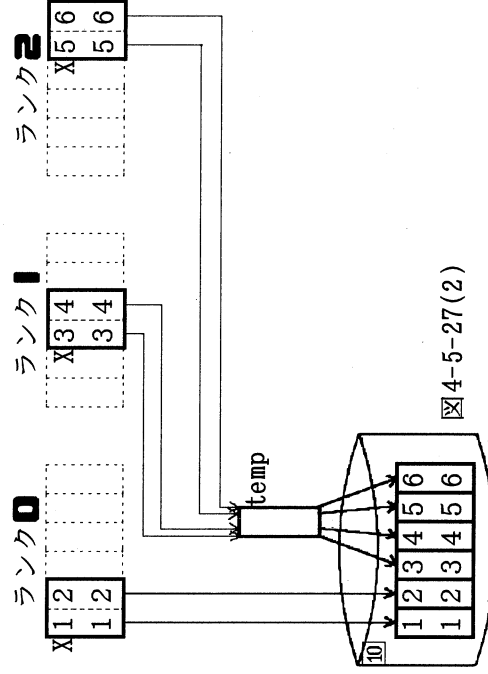


図4-5-27(2)

```

:
INCLUDE 'mpif.h'
REAL, ALLOCATABLE: TEMP(:, :)
REAL, ALLOCATABLE: Y(:, :)
PARAMETER (NCPU=3)
INTEGER JJSTA(0:NCPU-1), JJEND(0:NCPU-1)
:
DO IRANK=0, NPROCS-1
CALL PARA_RANGE(1,6,NPROCS, IRANK,
& JJSTA(IRANK), JJEND(IRANK))
ENDDO
JSTA = JJSTA(MYRANK)
JEND = JJEND(MYRANK)
ALLOCATE(Y(2, JSTA:JEND))
DO J=1, JSTA-1
READ(10) ← 空読み
ENDDO
DO J=JSTA, JEND
READ(10) (Y(I,J), I=1,2)
ENDDO
DO J=JSTA, JEND
DO I=1,2
~ = Y(I,J)
ENDDO
ENDDO
:

```

図4-5-28(1) 図4-5-28(2)のプログラム例

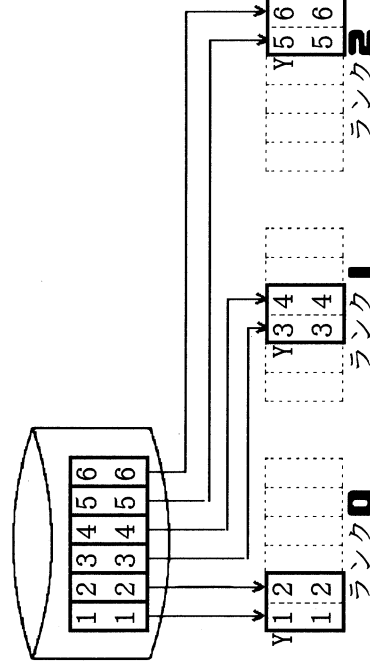


図4-5-28(2)

多重ループを並列化する場合、内側のループで並列化するか、外側のループで並列化するか、あるいは両方で並列化するかを決定する必要があります。本節ではこれを決定する際の考慮点について説明します。

■ 多重ループで内側と外側のどちらを並列化するか

まず以下の2つのD0ループを単体で実行したときの速度を比較します。図4-5-29(1)では内側のループが配列の左側の添字で反復しており、図4-5-29(2)ではその逆になっています。Fortranでは、2次元配列Aはメモリー上で A(1,1), A(2,1), ... の順に配置されますが、スカラー計算機では、図4-5-29(1)のようにメモリー上に配置された順番に要素を処理した方が、キャッシュ・ミスが発生しにくいため速くなります。一方C言語では、メモリー上で要素の配置順序が逆になるので、図4-5-29(2)の方が速くなります(キャッシュ・ミスの詳細は参考文献[6]の4章を参照)。

```

DIMENSION A(100,100)
DO J = 1, 100
  DO I = 1, 100
    A(I,J) = ~
  ENDDO
ENDDO
:

```

図4-5-29(1)

```

DIMENSION A(100,100)
DO I = 1, 100
  DO J = 1, 100
    A(I,J) = ~
  ENDDO
ENDDO
:

```

図4-5-29(2)

図4-5-29(1)をブロック分割で並列化するとして、内側と外側のどちらのループを並列化(計算量を1/nに縮小)するのがよいでしょうか？ 外側のJで並列化した例を図4-5-30(1)に、内側のIで並列化した例を図4-5-30(2)に示します。

配列を縮小しない場合を想定すると、前節の説明と重複しますが、図4-5-30(2)では各プロセスの②と③、④と⑤などがメモリー上で離れているため、計算および通信を行う際にキャッシュ・ミスが発生する可能性が高くなります。しかし図4-5-30(1)では①～⑥までの各要素がメモリー上で連続しているため、キャッシュ・ミスは発生しません。

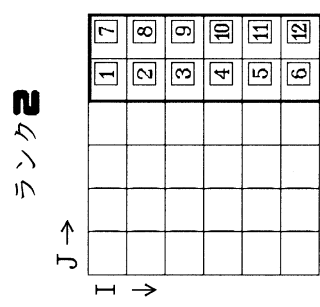
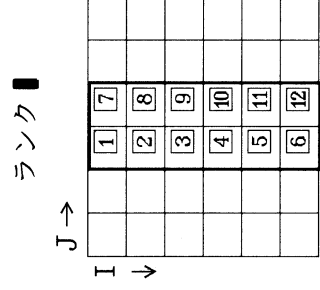
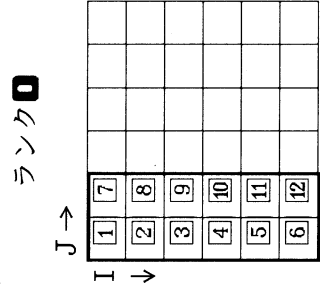
一方、配列を縮小している場合はこのような問題はありませんが、いずれにしても、他の条件が全く同じであれば、図4-5-30(1)のように外側のループを並列化するのがよいでしょう。

```

:
DIMENSION A(100,100)
DO J = JSTA, JEND
  DO I = 1, 100
    A(I,J) = ~
  ENDDO
ENDDO
:

```

図4-5-30(1)

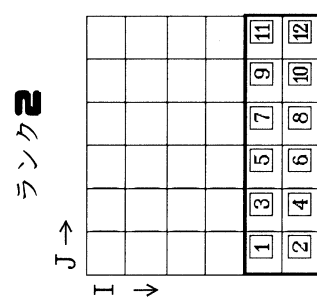
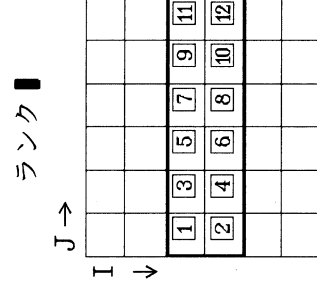
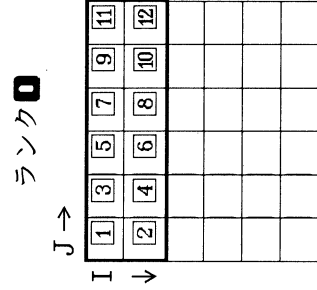


```

:
DIMENSION A(100,100)
DO J = 1, 100
  DO I = ISTA, IEND
    A(I,J) = ~
  ENDDO
ENDDO
:

```

図4-5-30(2)



■ 一方向にだけ通信が必要な場合

図4-5-31(1)では、図4-5-31(2)に示すように、配列Aの左と右の要素の値から配列Bの要素を計算します。このため、図4-5-31(3)のようにJで分割した場合、計算の前に、境界を越えた1列(黒の部分)を隣のプロセスから送信してもらう必要があります。これに対し、図4-5-31(4)のようにIで分割した場合、通信は必要ありません。このように、一方向にだけ通信が必要な場合には、通信が不要な方向に分割するようにします。このような例は差分法のプログラムで現れます。実際のプログラムでは、ある方向に通信が全く必要ないような場合はまずありますが、例えば計算領域が3次元でZ方向が重力方向の場合、プログラム内で並列化する全D0ループを検討してみると、XまたはY方向の通信に比べてZ方向の通信が少ない(あるいはその逆の)場合があります。

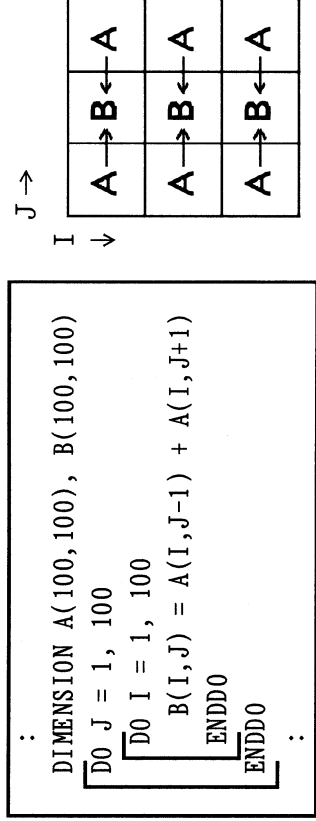


図4-5-31(1)

図4-5-31(2)

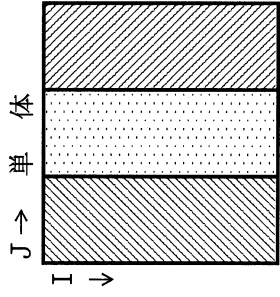


図4-5-31(3)

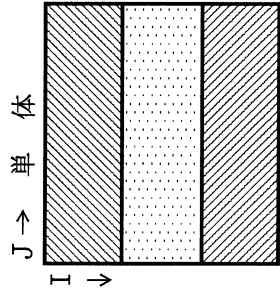
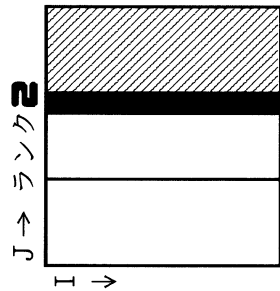
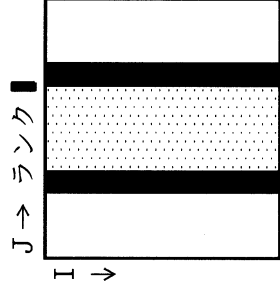
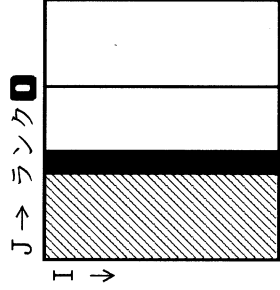
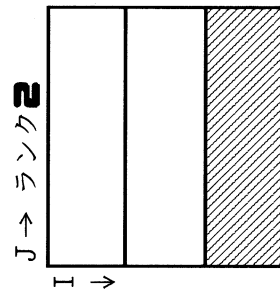
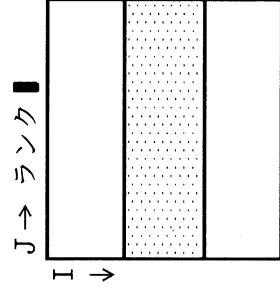
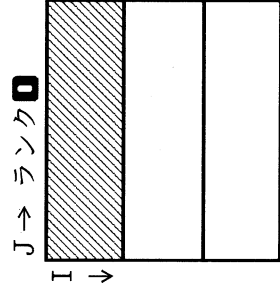


図4-5-31(4)



別の例として、例えば差分法の境界条件処理の部分が図4-5-31(5)のようになっていて、図4-5-31(6)のように境界の1列を反対側にコピーし、境界の行についてはコピーしないとします。この場合、図4-5-31(7)のように分割すると通信が発生しますが、図4-5-31(8)の分割では通信は発生しません。

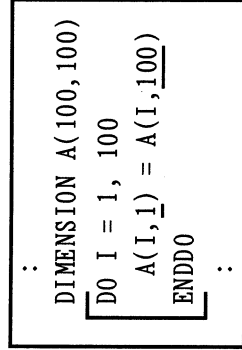


図4-5-31(5)

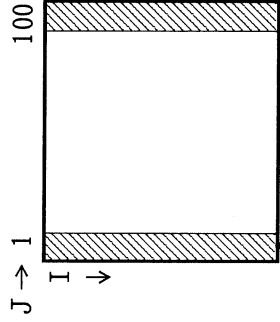


図4-5-31(6)

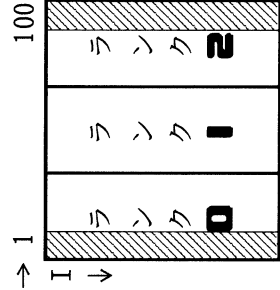


図4-5-31(7)

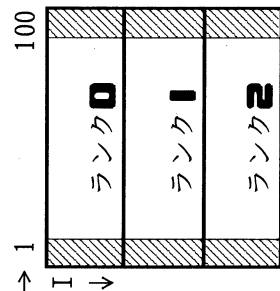


図4-5-31(8)

■ 計算領域が長方形(直方体)の場合

図4-5-32(1)では、図4-5-32(2)に示すように、配列Aの上下左右の値から配列Bを計算します。このため計算領域をIで分割してもJで分割しても通信が発生します。

ところで今までの例では、計算領域が正方形(3次元では立方体)の場合は扱ってきましたが、この例では計算領域が長方形(3次元では直方体)になっています(例えば新幹線のような細長い物体を解析する場合)。

I方向の方がJ方向よりも長いので、図4-5-32(3)のようにJで分割するよりも図4-5-32(4)のようにIで分割した方が、分割した境界の長さが短くなるため通信量も少なくなります。

このように、どちらの方向にも通信が必要な多重ループを分割する場合には、通信量が少なくなるような方向に分割するようにして下さい。

```

:
DIMENSION A(200,100), B(200,100)
DO J = 1, 100
  DO I = 1, 200
    B(I,J) = A(I-1,J) + A(I,J-1) + A(I,J+1) + A(I+1,J)
  ENDDO
ENDDO
:

```

図4-5-32(1)

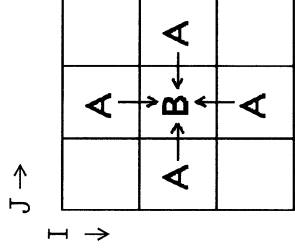


図4-5-32(2)

単体

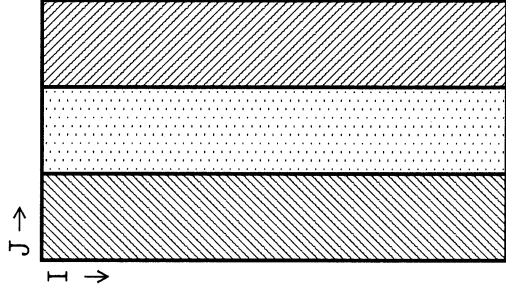
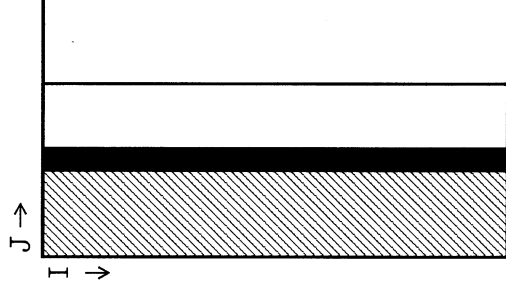
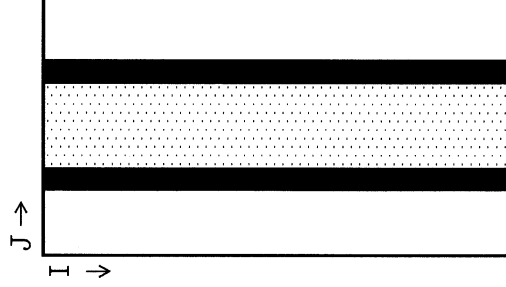


図4-5-32(3)

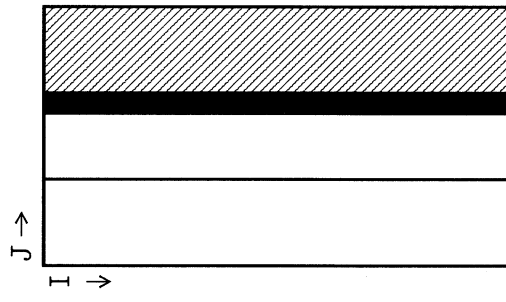
ランク0



ランク1



ランク2



単体

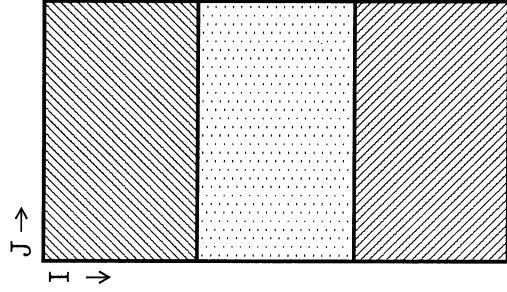
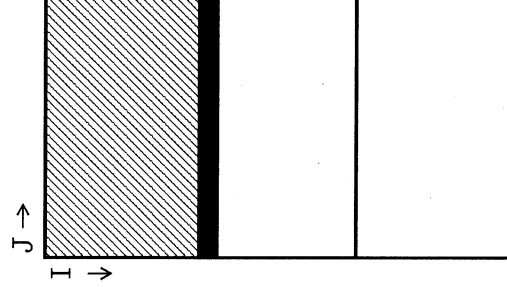
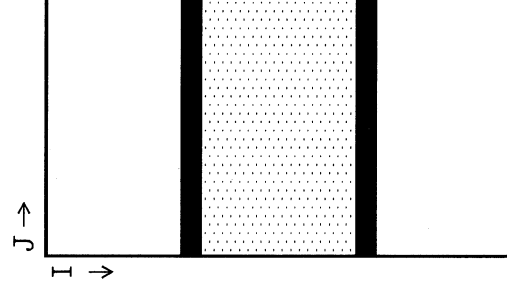


図4-5-32(4)

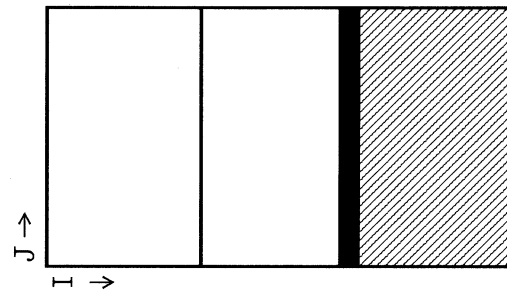
ランク0



ランク1



ランク2



## ■ 2方向のブロック分割

今までは全て、IまたはJのどちらか1方向でブロック分割を行いましたが、4-5-2節で紹介したようにIとJの2方向でブロック分割を行うこともできます。プログラム例を図4-5-33(1)に示します。

1方向だけで分割した図4-5-33(2)では、プロセッサ数を増やしていても境界の長さは変わらないので、次第に通信の割合が増えてサチュレートします。一方2方向でブロック分割した図4-5-33(3)では、プロセッサ数が増えるにつれて境界の長さが短くなり通信量が減るため、サチュレートしにくくなります(ただし通信回数は2倍になって1方向だけの分割より時間がかかり、また修正も面倒になります)。従って、差分法でロード数が多い場合に良好なパフォーマンスを得るためには、2方向の分割が必須です。なお、この分割はロード数が整数×整数で表せる場合(例えば $4 \times 2 = 8$ )にのみ可能です。適用例については5-1-3節を参照して下さい。

```

DO J = JSTA, JEND
DO I = ISTA, IEND
  B(I,J) = A(I-1,J) + A(I,J-1) + A(I,J+1) + A(I+1,J)
ENDDO
ENDDO

```

図4-5-33(1)

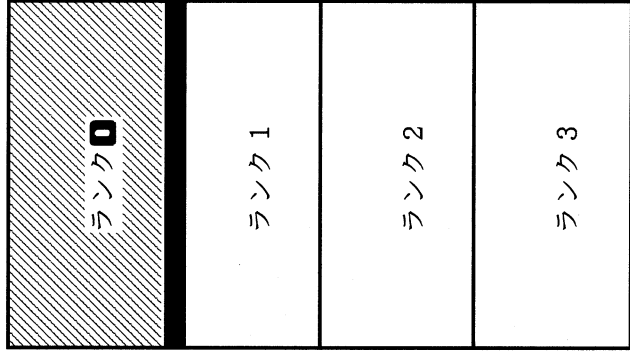


図4-5-33(2)

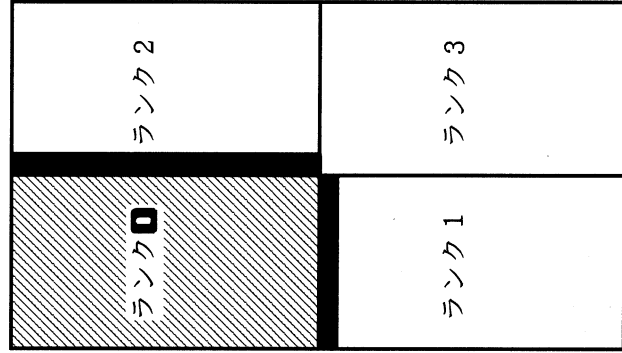
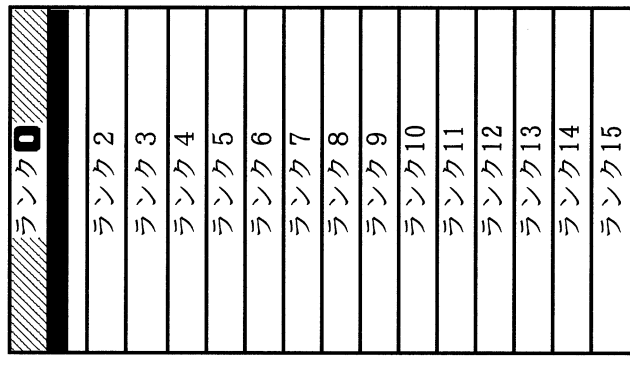
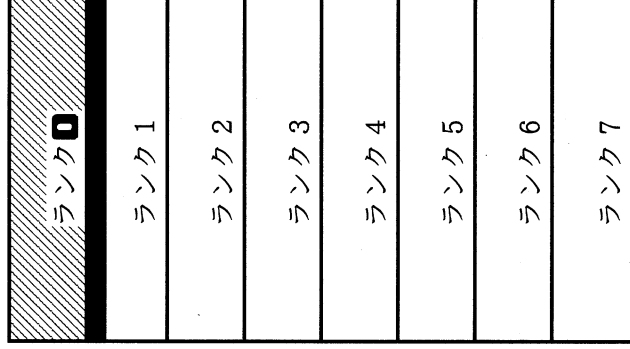
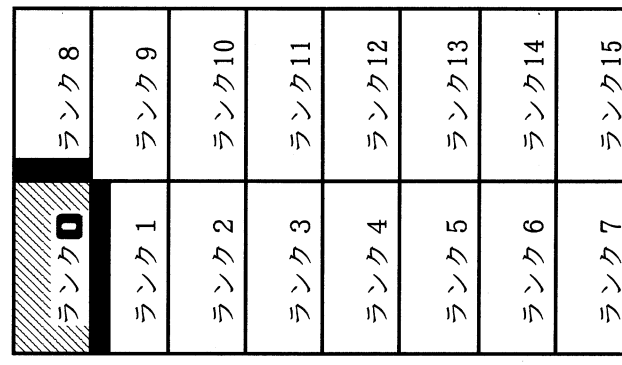
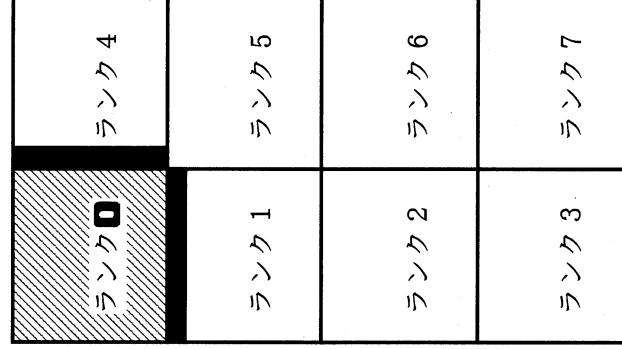


図4-5-33(3)



前節では、本当の意味の並列化の方法、すなわち計算量を $1/n$ に減少させる方法について説明し、その具体的な方法としてブロック分割、サイクリック分割などの方法を紹介しました。一方、1章で説明しましたが、『本当の意味の並列化に伴う矛盾(副作用)を解消するために、必要最小限、仕方なしに行う』のがメッセージ交換であり、メッセージ交換サブルーチンの使い方は3章で説明しました。本節では、メッセージ交換が実際のプログラムのどのような局面で必要になり、具体的にどのようなように行うのかを説明します。

### 4-6-1 メッセージ交換はどんなときに必要か

メッセージ交換がどんなときに必要となるかについて、原点に戻って考えてみましょう。図4-6-1(1)の単体プログラムでは、D0ループで配列Aの値を更新し、D0ループから抜けた後の①でA(6)の値を参照します。これを図4-6-1(4)のように表現します。

このプログラムを並列化すると図4-6-1(2)のようになります。①のステートメントは全プロセスが実行しますが、図4-6-1(5)に示すように、A(6)に正しい値が入っているのはランク2のプロセスのみなので、このままでは(ランク0とランク1のプロセスでは)結果がおかしくなってしまいます。これが1章で説明した「本当の意味の並列化に伴う矛盾(副作用)」です。

そこで図4-6-1(3)に示すように、並列化したD0ループの終了後の②で、ランク2のプロセスがA(6)をランク0と1のプロセスに送信し、その後で各プロセスが①を実行すれば正しく動作します。これが「本当の意味の並列化に伴う矛盾(副作用)を解消するために、必要最小限、仕方なしに行うメッセージ交換」です。

以上をまとめます。D0ループに対して本当の意味の並列化すると、各プロセスでは計算を $1/n$ (ノード数がn台の場合)しか行わないので、そのD0ループより後の部分で矛盾(副作用)が発生する可能性があり、それを解消するためにメッセージ交換を行います。そして、プログラムのどの部分で、どの変数を、どのプロセス間で通信すれば、本当の意味の並列化に伴う矛盾(副作用)を解消できて、しかも通信量を最小にできるか、がポイントであり、これを決定するのはユーザーということになります。

このように一般的に説明すると、メッセージ交換にはいろいろなバリエーションがあって難しく感じられるかも知れませんが、普通のプログラムではそれほど多くのパターンがあるわけではありません。典型的なメッセージ交換のパターンについて次節以降で説明します。

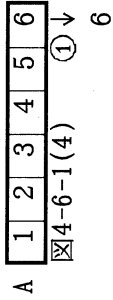
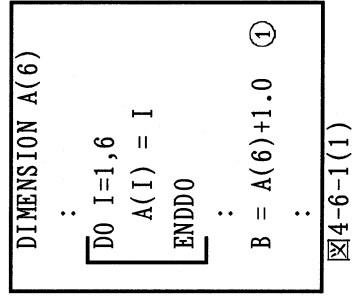


図4-6-1(4)

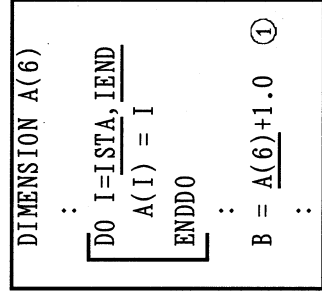


図4-6-1(2)

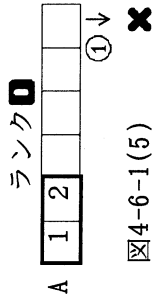


図4-6-1(5)

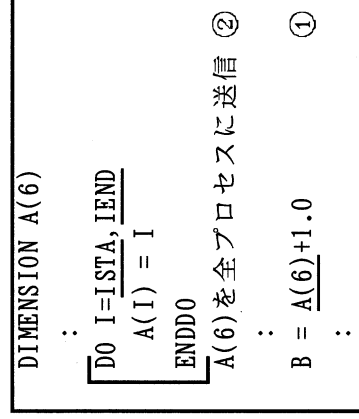
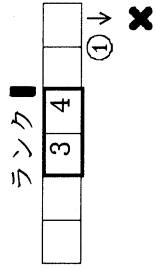


図4-6-1(3)

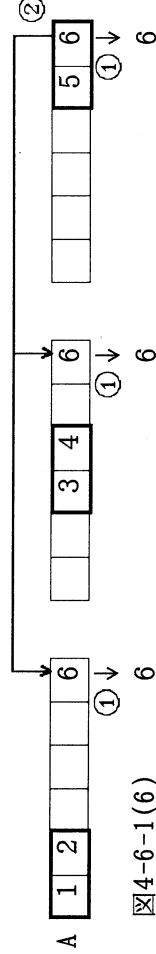
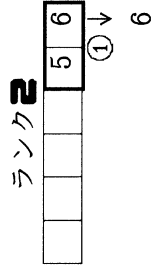


図4-6-1(6)

## 4-6-2 1次元差分法

キングオプ科学技術計算の称号を差分法に与えても異議をとない人はいないでしょう。実際、過去に私が並列化したプログラムのうち、約半数が差分法のプログラムでした。

4-3節の「計算パターン2」で簡単に説明したように、差分法を並列化する場合、境界のデータを隣接するプロセス間でお互いに交換するパターン(以後シフト)が現れます。本節では最も簡単な1次元の差分法でのシフトについて説明し、一般的な2次元のシフトについては5-1節で説明します。

### ■ シフト(1次元の場合)

差分法モドキの簡単なプログラムとして図4-6-2(1)を例に取り上げます。まず①で配列Aに適当な値を設定し、次に②で、配列Aの前後の値から配列Bの値を求めます。データの動きを図示すると図4-6-2(2)のようになります。

これを例えば3つのプロセスで並列に計算する場合、配列AとBを例えば図4-6-2(3)のようにブロック分割し、各プロセスは自分の担当する要素のみを計算します。例えばランク■は、①ではA(5)~A(8)を計算し、②ではB(5)~B(8)を計算します。②の計算でA(4)~A(9)の値を参照しますが、A(4)は①でランク■が計算し、A(9)は①でランク■が計算しているので、ランク■はこれらの値を持っていません。そこで②の計算に入る前に、A(4)をランク■から、A(9)をランク■からそれぞれ送ってもらわなければならない。逆に、ランク■はA(5)をランク■に、A(8)をランク■に送信してあげる必要があります。このように、②の計算に入る前に境界のデータを隣接するプロセス間でお互いにシフトする必要があります。

なお、この例では、例えばランク■のプロセスはA(4)とA(9)の値を①で計算しておけばシフトの必要はありませんが、実際のプログラムではほとんどの場合シフトが必要になると考えて下さい。

```

PROGRAM MAIN
  IMPLICIT REAL*8(A-H,O-Z)
  PARAMETER (N=11)
  DIMENSION A(N), B(N)
  DO I = 1, N
    A(I) = I
  ENDDO
  DO I = 2, N-1
    B(I) = A(I-1) + A(I+1)
  ENDDO
END
  
```

図4-6-2(1)

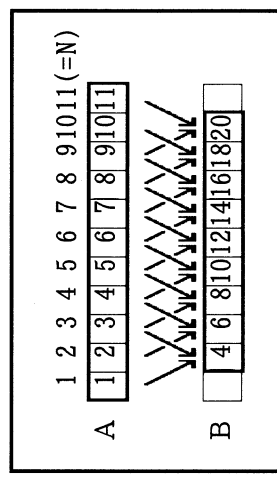


図4-6-2(2)

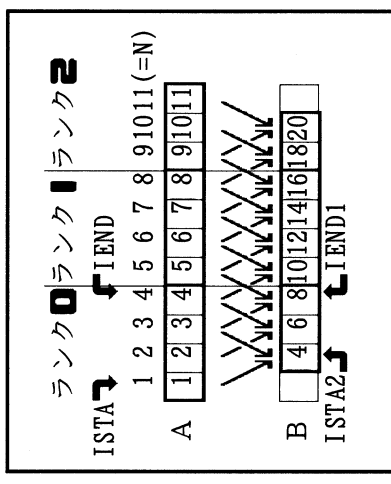


図4-6-2(3)

図4-6-2(1)を並列化したプログラムを図4-6-3(1)に、データの動きを図4-6-3(4)に示します(図4-6-3(1)と図4-6-3(4)の⑥~⑨は対応しています)。

●まず各プロセスが担当する配列AとBの範囲を決定します。このとき図4-6-2(3)に示したように、配列AとBでは分割した境界の位置(図中の縦線部分)が一致していなければなりません。しかし配列A(①のループ)の範囲は1~N、配列B(②のループ)の範囲は2~N-1と異なっているため、配列Aを基準にして分割し、それに配列Bの分割の境界を合わせるか、その逆にするか、いずれかを選択します(どちらでも構いません)。ここでは前者の方法を選択することにします。



図4-6-3(1)の③で、まずサブルーチンPARA\_RANGE(4-5-4節参照)を使用して配列Aの範囲(1~N)をブロック分割し、各プロセスが担当する配列A(①)の反復の下限(ISTA)と上限(IEND)を取得します。③の残りの部分では各プロセスが担当する配列B(②)のループの反復の下限(ISTA2)と上限(IEND1)を決定します。図4-6-2(3)の縦線を一致させるため、端以外のプロセスではISTA2をISTAと、IEND2をIENDと等しくし、端を担当するランク**0**のプロセスのISTA2とISTA、およびランク**2**のプロセスのIEND1とIENDだけが異なるようにします。例えばランク**0**のプロセスの担当範囲ISTA, IEND, ISTA2, IEND1は図4-6-2(3)のようになります。なお、この例では配列Aを基準に分割し配列Bの分割をそれに合わせたため、配列Bの分割のロードバランズが各プロセス間で不均等になっていますが、要素数(N)が多くなればロードバランズはほぼ均等になり、問題はありません。

● ③で得たISTA, IEND, ISTA2, IEND1を使用して、図4-6-2(1)の①と②のループ反復を、図4-6-3(1)の⑤と⑩のように修正します。

● 前述のように、⑩の計算に入る前に、境界のデータを隣接するプロセス間でシフトする必要があります。シフトで通信されるデータの動きを図4-6-3(4)の(実線の)『 $\square$ 』と『 $\sqcup$ 』に示します。

以後『 $\square$ 』のシフトについて説明します。まず図4-6-3(1)の④で、『 $\square$ 』方向の宛先プロセスのランク値を変数IUPに、送信元プロセスのランク値を変数IDOWNに設定します。ところでこの例では『 $\square$ 』の点線部分のシフトは行なっていません。これを非循環シフトと呼び、点線部分のシフトを行う場合を循環シフトと呼びます。多くのプログラムでは非循環シフトになっています。非循環シフトでは、一番右端(ランクNPROCS-1)のプロセスは『 $\square$ 』の送信をしないので、IUPにはMPI\_PROC\_NULLを指定します(付録のMPI\_ISEND, MPI\_IRECVの引数参照)。同様に一番左端(ランク**0**)のプロセスは『 $\square$ 』の受信をしないので、IDOWNにはMPI\_PROC\_NULLを指定します。

● 図4-6-3(1)の⑥~⑩の二重線で囲んだ部分でシフトを行います。なお、シフトには以下で説明する通信方法以外に、(注)で説明する2つの方法があり、どの方法が最も速いかはマシンの環境によって異なります。

⑥では『 $\square$ 』の送信を行います。宛先プロセスのランクはIUP、送信するデータは図4-6-3(4)に示すようにA(IEND)となります。この送信に対応する『 $\square$ 』の受信は⑧で、送信元プロセスのランクはIDOWN、受信したデータを入れる場所は図4-6-3(4)に示すようにA(ISTA-1)となります。

同様に⑦と⑨で『 $\sqcup$ 』の送信と受信を行います。宛先プロセスと送信元プロセスが『 $\square$ 』の場合と逆になるので注意して下さい。

⑥~⑨は非ブロッキング通信ルーチンなので一気に実行(が開始)されます。そしてそれぞれに対し、ストッパーとして⑩でMPI\_WAITを実行します。MPI\_WAITの引数で配列ISTATUSを使用するので、ISTATUSを配列宣言するのを忘れないで下さい。

● ⑥~⑩で⑩の計算に必要なデータがそろったので⑩を計算します。

以上で図4-6-3(1)の説明を終了します。このように、差分法で並列化するD0ループ内のステートメントの右辺にA(I-1)やA(I+1)のような項がある場合、そのループに入る前に境界のデータをシフトするのが基本的なパターンとなります。言うまでもありませんが、例えばA(I-1)のみが現れている場合は図4-6-3(4)の『 $\sqcup$ 』のシフト、すなわち⑥、⑧とMPI\_WAITのみを行います。

なお、本例は1次元なのでプロセス間の境界の1要素(点)をシフトしましたが、2次元の場合は『線』を、3次元の場合は『面』をシフトすることになります。

### ■ 別の通信方法(重要)

図4-6-3(1)の二重線で囲んだ部分の通信(シフト)を行う別の方法を紹介します。まず図4-6-3(1)の⑥~⑩の順序を変更した図4-6-3(2)で、同じ通信を行うことができます。

またMPIではシフト専用の1対1通信サブルーチンMPI\_SENDRECV(付録参照)が提供されており、このルーチンを使用した場合は図4-6-3(3)のようになります。例えば図4-6-3(3)の□は図4-6-3(2)の□と同じです。このルーチンはブロッキング通信ルーチンですが、デッドロックにはなりません。

以上3つの通信方法のうち、どの方法がもっとも速いかは、マシンの環境によって異なります。なお、本書の以後のプログラム例では、シフトの通信は便宜的に図4-6-3(1)で行いました。

```

PROGRAM MAIN
IMPLICIT REAL*8(A-H,0-Z)
INCLUDE 'mpif.h'
PARAMETER (N=11)
DIMENSION A(N),B(N)
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(~, NPROCS, IERR)
CALL MPI_COMM_RANK(~, MYRANK, IERR)
CALL PARA_RANGE(1, N, NPROCS, MYRANK, ISTA, IEND)
ISTA2 = ISTA
IEND1 = IEND
IF (MYRANK == 0 ) ISTA2 = 2
IF (MYRANK == NPROCS-1) IEND1 = N-1
IUP = MYRANK + 1
IDOWN = MYRANK - 1
IF (MYRANK == NPROCS-1) IUP = MPI_PROC_NULL
IF (MYRANK == 0) IDOWN = MPI_PROC_NULL
DO I = ISTA, IEND
  A(I) = I
ENDDO
CALL MPI_ISEND(A(IEND), 1, MPI_REAL8,
& IUP, 1, MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(A(ISTA), 1, MPI_REAL8,
& IDOWN, 1, MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_IRECV(A(ISTA-1), 1, MPI_REAL8,
& IDOWN, 1, MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_IRECV(A(IEND+1), 1, MPI_REAL8,
& IUP, 1, MPI_COMM_WORLD, IRECV2, IERR)
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
DO I = ISTA2, IEND1
  B(I) = A(I-1) + A(I+1)
ENDDO
CALL MPI_FINALIZE(IERR)
END

```

図4-6-3(1) 通信方法1

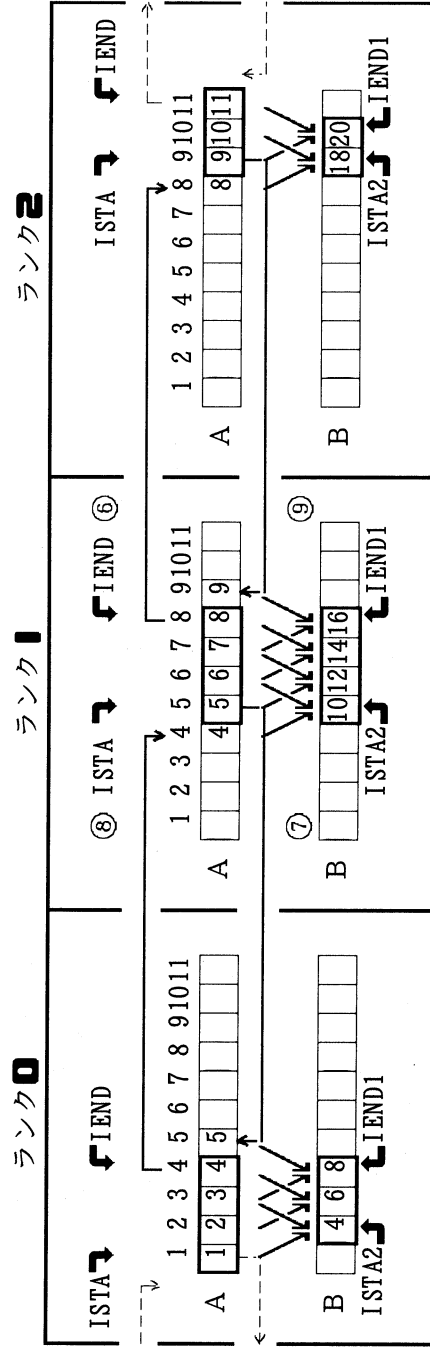


図4-6-3(4)

```

:
CALL MPI_ISEND(A(IEND), 1, MPI_REAL8,
& IUP, 1, MPI_COMM_WORLD, ISEND, IERR)
CALL MPI_IRECV(A(ISTA-1), 1, MPI_REAL8,
& IDOWN, 1, MPI_COMM_WORLD, IRECV, IERR)
CALL MPI_WAIT(ISEND, ISTATUS, IERR)
CALL MPI_WAIT(IRECV, ISTATUS, IERR)
CALL MPI_ISEND(A(ISTA), 1, MPI_REAL8,
& IDOWN, 1, MPI_COMM_WORLD, ISEND, IERR)
CALL MPI_IRECV(A(IEND+1), 1, MPI_REAL8,
& IUP, 1, MPI_COMM_WORLD, IRECV, IERR)
CALL MPI_WAIT(ISEND, ISTATUS, IERR)
CALL MPI_WAIT(IRECV, ISTATUS, IERR)
:

```

図4-6-3(2) 通信方法2

```

:
CALL MPI_SENDREC
& (A(IEND), 1, MPI_REAL8, IUP, 1,
& A(ISTA-1), 1, MPI_REAL8, IDOWN, 1,
& MPI_COMM_WORLD, ISTATUS, IERR)
CALL MPI_SENDREC
& (A(ISTA), 1, MPI_REAL8, IDOWN, 1,
& A(IEND+1), 1, MPI_REAL8, IUP, 1,
& MPI_COMM_WORLD, ISTATUS, IERR)
:

```

図4-6-3(3) 通信方法3

■ 境界条件の処理

差分法を並列化した場合、境界条件の処理を行う部分で以下のような修正が必要になることがあります。まず図4-6-4(1)の①では、境界の要素A(1)だけが特別な処理をしています。並列化した図4-6-4(2)でA(1)を担当しているのはランク0のプロセスなので、②のIF文を追加して、ランク0のプロセスのみが①の処理を行うようにします。

```

:
DO I = 1, N
  A(I) = B(I) + C(I)
ENDDO
A(1) = A(1) + 1.0 ①
:
    
```

図4-6-4(1)

```

:
DO I = ISTAR, IEND
  A(I) = B(I) + C(I)
ENDDO
IF (MYRANK==0) A(1) = A(1) + 1.0 ②
:
    
```

図4-6-4(2)

次に、境界条件が周期的境界になっている場合の例を示します。図4-6-5(1)と図4-6-6(1)の①では、各境界のデータを反対側の境界にコピーしています。これを並列化した場合、図4-6-5(2)に示すように、A(0), A(1)を担当するランク0のプロセスと、A(N), A(N+1)を担当するランク2のプロセスの間で、 $\hookrightarrow$ と $\leftarrow$ の通信が必要となります。この通信を行うプログラムを図4-6-6(2)に示します。なお、このプログラムを2プロセス以上で実行した場合は問題ありませんが、1プロセスで実行すると誤動作します。参考までにその理由と対処方法を以下に示します。いずれにしても、並列化したプログラムを1プロセスで実行すると、他の箇所でも同様の問題が発生する可能性があるため、3-4-5節の「1対1通信を1プロセスで実行した場合」参照。

図4-6-6(2)を1プロセスで実行すると、④がELSEIFになっていてランク0は⑤と⑥の通信を実行しないので、②から③への通信のみがランク0内で行われてしまい、結果は図4-6-5(3)のようになります。そこで、ELSEIFを分離して別のIF文にすると図4-6-6(3)のようになります。ここで本来は⑦→⑩と⑩→⑧の通信が行われるはずですが、⑦, ⑧, ⑩, ⑪の下線に示すタグの値が同一(「1」)のため、⑦→⑧と⑩→⑪の通信が行われてしまい、結果は図4-6-5(4)のようになってしまいます。

そこで、図4-6-6(4)の⑫, ⑬, ⑭, ⑮の下線に示すように、送信と受信のペアごとにタグの値を区別して見ます。プログラムを実行すると、⑬の受信命令が実行された後、その受信が実際に行われるまで⑭でWAITします。このWAITは⑮の送信が行われれば解除されますが、⑭でWAITしているため⑮が実行出来ず、デッドロックになってしまいます(これは3-4-4節の「パターン2」に相当します)。

そこで、図4-6-6(5)のように、送信命令と受信命令を全て実行してからWAITするようになれば、デッドロックを回避でき、1プロセスでも正常に動作します。境界を担当する2プロセス間で通信が必要となる並列プログラムを、1プロセスで実行する可能性がある場合、図4-6-6(5)のようにプログラムが複雑になってしまいます。このような場合は、図4-6-6(6)のように、プロセス数が2プロセス以上の場合と1プロセスの場合で処理を分けた方が良いかも知れません。

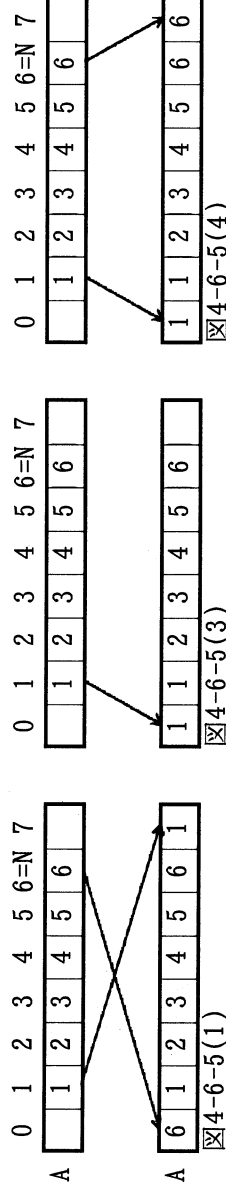
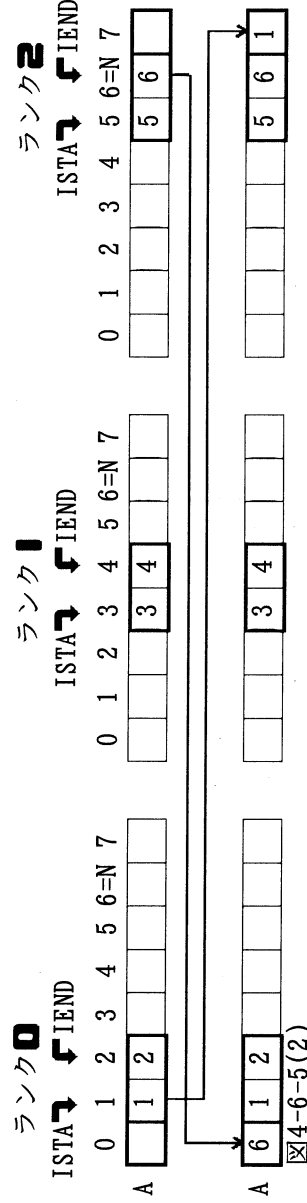


図4-6-5(1)

図4-6-5(2)

図4-6-5(4)



ランク0

ランク2

```

:
PARAMETER(N=6)
REAL A(0:N+1)
DO I=1,N
  A(I) = ~
ENDDO
A(0) = A(N)
A(N+1) = A(1)
:

```

☒4-6-6(1)

```

:
IF (MYRANK==0) THEN
CALL MPI_ISEND(A(1),1,MPI_REAL,NPROCS-1,1,
& MPI_COMM_WORLD,IREQ1,IERR) ②
CALL MPI_IRECV(A(0),1,MPI_REAL,NPROCS-1,1,
& MPI_COMM_WORLD,IREQ2,IERR) ③
CALL MPI_WAIT(IREQ1,ISTATUS,IERR)
CALL MPI_WAIT(IREQ2,ISTATUS,IERR)
ELSEIF (MYRANK==NPROCS-1) THEN ④
CALL MPI_ISEND(A(N),1,MPI_REAL,0,1,
& MPI_COMM_WORLD,IREQ1,IERR) ⑤
CALL MPI_IRECV(A(N+1),1,MPI_REAL,0,1,
& MPI_COMM_WORLD,IREQ2,IERR) ⑥
CALL MPI_WAIT(IREQ1,ISTATUS,IERR)
CALL MPI_WAIT(IREQ2,ISTATUS,IERR)
ENDIF
:

```

☒4-6-6(2)

```

:
IF (MYRANK==0) THEN
CALL MPI_ISEND(A(1),1,MPI_REAL,NPROCS-1,1,
& MPI_COMM_WORLD,IREQ1,IERR) ⑦
CALL MPI_IRECV(A(0),1,MPI_REAL,NPROCS-1,1,
& MPI_COMM_WORLD,IREQ2,IERR) ⑧
CALL MPI_WAIT(IREQ1,ISTATUS,IERR)
CALL MPI_WAIT(IREQ2,ISTATUS,IERR)
ENDIF
IF (MYRANK==NPROCS-1) THEN ⑨
CALL MPI_ISEND(A(N),1,MPI_REAL,0,1,
& MPI_COMM_WORLD,IREQ1,IERR) ⑩
CALL MPI_IRECV(A(N+1),1,MPI_REAL,0,1,
& MPI_COMM_WORLD,IREQ2,IERR) ⑪
CALL MPI_WAIT(IREQ1,ISTATUS,IERR)
CALL MPI_WAIT(IREQ2,ISTATUS,IERR)
ENDIF
:

```

☒4-6-6(3)

```

:
IF (MYRANK==0) THEN
CALL MPI_ISEND(A(1),1,MPI_REAL,NPROCS-1,1,
& MPI_COMM_WORLD,IREQ1,IERR) ⑫
CALL MPI_IRECV(A(0),1,MPI_REAL,NPROCS-1,2,
& MPI_COMM_WORLD,IREQ2,IERR) ⑬
CALL MPI_WAIT(IREQ1,ISTATUS,IERR)
CALL MPI_WAIT(IREQ2,ISTATUS,IERR) ⑭
ENDIF
IF (MYRANK==NPROCS-1) THEN
CALL MPI_ISEND(A(N),1,MPI_REAL,0,2,
& MPI_COMM_WORLD,IREQ1,IERR) ⑮
CALL MPI_IRECV(A(N+1),1,MPI_REAL,0,1,
& MPI_COMM_WORLD,IREQ2,IERR) ⑯
CALL MPI_WAIT(IREQ1,ISTATUS,IERR)
CALL MPI_WAIT(IREQ2,ISTATUS,IERR)
ENDIF
:

```

☒4-6-6(4)

```

:
IF (MYRANK==0) THEN
CALL MPI_ISEND(A(1),1,MPI_REAL,NPROCS-1,1,
& MPI_COMM_WORLD,IREQ1,IERR)
CALL MPI_IRECV(A(0),1,MPI_REAL,NPROCS-1,2,
& MPI_COMM_WORLD,IREQ2,IERR)
ENDIF
IF (MYRANK==NPROCS-1) THEN
CALL MPI_ISEND(A(N),1,MPI_REAL,0,2,
& MPI_COMM_WORLD,IREQ3,IERR)
CALL MPI_IRECV(A(N+1),1,MPI_REAL,0,1,
& MPI_COMM_WORLD,IREQ4,IERR)
ENDIF
IF (MYRANK==0) THEN
CALL MPI_WAIT(IREQ1,ISTATUS,IERR)
CALL MPI_WAIT(IREQ2,ISTATUS,IERR)
ENDIF
IF (MYRANK==NPROCS-1) THEN
CALL MPI_WAIT(IREQ3,ISTATUS,IERR)
CALL MPI_WAIT(IREQ4,ISTATUS,IERR)
ENDIF
:

```

☒4-6-6(5)

```

:
IF (NPROCS==1) THEN
A(0) = A(N)
A(N+1) = A(1)
ELSE
☒4-6-6(2)と同じ
ENDIF
:

```

☒4-6-6(6)

### 4-6-3 集団通信サブルーチンの代替

#### 4-6-3-1 MPI\_GATHERVの代替

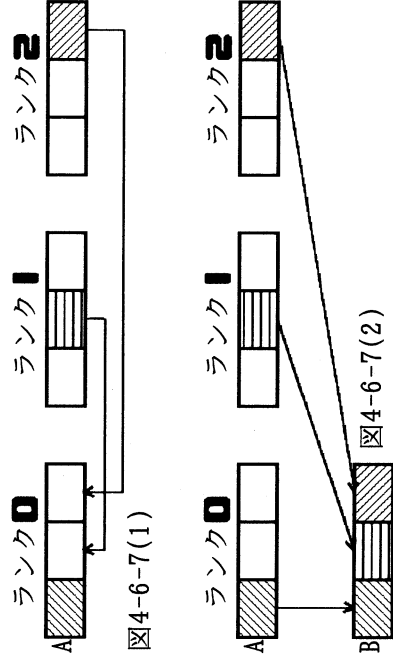
各プロセスが計算した結果を、例えばランク0に収集する通信パターンが現れます(4-4-2節参照)。

#### ■ (ケース1) MPI-2のMPI\_IN\_PLACEが使用でき、データが連続な場合

配列を縮小していない場合、図4-6-7(1)のように、ランク0の配列Aの空いている部分に収集できればよいのですが、集団通信サブルーチンには『送信バッファと受信バッファがメモリー上で重なってはいけない』(3-3-5節参照)という制限があるため、図4-6-7(2)のように、新たに受信バッファ用の配列Bを設ける必要があります。しかし、MPI-2のMPI\_IN\_PLACE(3-8-1節参照)という変数を使用すると、この制限を回避することができます(MPI-2が使用できない環境では使用できません)。

図4-6-7(3)に示すM×Nの2次元配列Aを2次元方向に分割する場合は想定します。一般にプロセスによってデータ量が異なるため、MPI\_GATHERではなくMPI\_GATHERVを使用して、各プロセスのデータをランク0に集めます。図4-6-7(4)の[4]で以下の2つの配列を作成し、ランク0は[2]で、その他のプロセスは[3]で、MPI\_GATHERV(付録の「MPI\_IN\_PLACEの使い方」参照)をコールします。[2]では送信バッファの代わりにMPI\_IN\_PLACEを指定します。[3]では受信バッファは不要なので、[3]の下線部は実際には無視されます。

- IRCNT(i)：ランクiのプロセスが担当する要素数を設定します。
- IDISP(i)：ランクiのプロセスが担当する最初の要素の、配列の先頭からの変位(回,⑧,⑭)を設定します。



```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER (M=~, N=~, NCPU=3)
  REAL A(M,N)

  INTEGER IRCNT(0:NCPU-1), IDISP(0:NCPU-1) [1]
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~, NPROCS, IERR)
  CALL MPI_COMM_RANK(~, MYRANK, IERR)
  DO IRANK=0, NPROCS-1
    CALL PARA_RANGE
    & (1, N, NPROCS, IRANK, JSTA, JEND)
    IRCNT(IRANK) = M*(JEND-JSTA+1) [1]
    IDISP(IRANK) = M*(JSTA-1)
  ENDDO
  CALL PARA_RANGE
  & (1, N, NPROCS, MYRANK, JSTA, JEND)
  :
  IF (MYRANK==0) THEN
    CALL MPI_GATHERV
    & (MPI_IN_PLACE, IDUMMY, IDUMMY, [2]
    & A, IRCNT, IDISP, MPI_REAL,
    & 0, MPI_COMM_WORLD, IERR)
  ELSE
    CALL MPI_GATHERV
    & (A(1, JSTA), IRCNT(MYRANK), MPI_REAL,
    & IDUMMY, IDUMMY, IDUMMY, [3]
    & 0, MPI_COMM_WORLD, IERR)
  ENDIF
  :
  
```

図4-6-7(4)

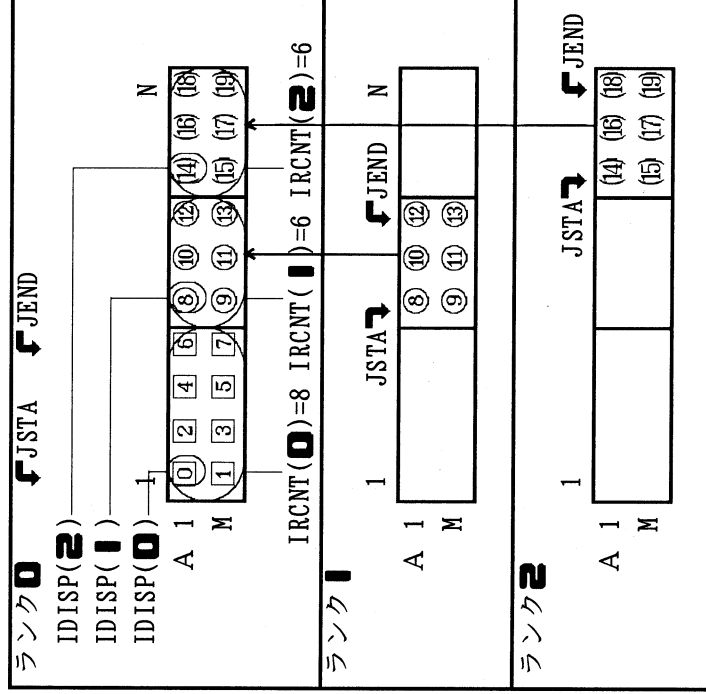


図4-6-7(3)

■ (ケース2) MPI-2のMPI\_IN\_PLACEが使用できず、データが連続な場合

MPI-2のMPI\_IN\_PLACEの機能が使用できないマシン環境で、前述のケース1のように、配列Bを新たに設けず配列Aの空いている部分に直接収集したい場合、MPI\_GATHERVと同等の通信を、1対1通信サブルーチンを使用して行います。

まず図4-6-8(2)の[1]で以下の配列を作成します。

- JJSTA(i)：ランクiのプロセスが担当する部分の、2次元方向の最初の要素番号を設定します。
- JJLEN(i)：ランクiのプロセスが担当する要素数を設定します。

[2]と[3]がMPI\_GATHERVの代替で通信を行う部分です。図4-6-8(1)からわかるように、ランク**0**以外のプロセスは、[3]で自分が担当したデータをランク**0**のプロセスに送信し、ランク**0**のプロセスは、[2]で自分以外のプロセスから送られたデータを配列Aのしかるべき部分に受信します。[2]の下線部に示すように、D0ルーブが「0」からではなく「1」から開始することに注意して下さい。

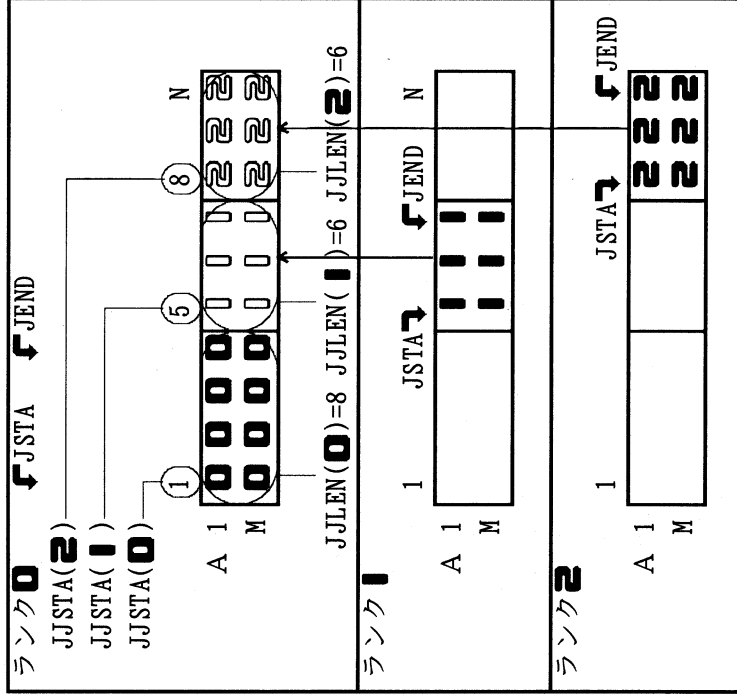


図4-6-8(1)

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER(M=~,N=~,NCPU=3)
  REAL A(M,N)
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  INTEGER JJSTA(0:NCPU-1),JJLEN(0:NCPU-1) [1]
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~,NPROCS,IERR)
  CALL MPI_COMM_RANK(~,MYRANK,IERR)
  DO IRANK=0,NPROCS-1
    CALL PARA_RANGE
    & (1,N,NPROCS,IRANK,JSTA,JEND)
    JJSTA(IRANK) = JSTA
    JJLEN(IRANK) = M*(JEND-JSTA+1) [1]
  ENDDO
  CALL PARA_RANGE
  & (1,N,NPROCS,MYRANK,JSTA,JEND)
  :
  IF (MYRANK==0) THEN
    DO IRANK=1,NPROCS-1
      CALL MPI_RECV(A(1,JJSTA(IRANK)),
        JJLEN(IRANK),MPI_REAL,IRANK,1,
        MPI_COMM_WORLD,IREQ,IERR) [2]
      CALL MPI_WAIT(IREQ,ISTATUS,IERR)
    ENDDO
  ELSE
    CALL MPI_SEND(A(1,JSTA),
      JJLEN(MYRANK),MPI_REAL,0,1,
      MPI_COMM_WORLD,IREQ,IERR) [3]
    CALL MPI_WAIT(IREQ,ISTATUS,IERR)
  ENDIF
  :
  
```

図4-6-8(2)

■ (ケース3) データが不連続で派生データ型を使う場合

ケース1とケース2では、各プロセスの送信データがメモリー上で連続してしました。図4-6-9(1)では配列Aを1次元方向に分割しているので、各プロセスの送信データはメモリー上でとびとびになっています (Fortranの場合)。とびとびのデータを、いったん連続した一時配列に入れてから送信(受信はその逆)する方法もありますが、それが面倒な場合、3-5節で説明した派生データ型を使用します。

図4-6-9(3)の[4]で、MPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照)(MPI-2が使用できる場合のみ)を使用して、ランクiのプロセスの担当部分の派生データ型を配列ITYPE(i)に作成します(図4-6-9(2)参照)。ケース1で使用したMPI\_GATHERVでは、受信バッファのデータ型を1種類しか指定できない(つまりプロセスごとに配列で別々に指定できない)ため、ケース2と同じ方法で通信を行います。

[2]のMPI\_IRECVは、『ランクIRANKのプロセスから送られたデータを、配列A(の先頭)から始まる1個のITYPE(IRANK)型の受信バッファに受信する』という意味になり、[3]のMPI\_ISENDは、『配列A(の先頭)から始まる1個のITYPE(MYRANK)型のデータを、ランク0のプロセスに送信する』という意味になります。

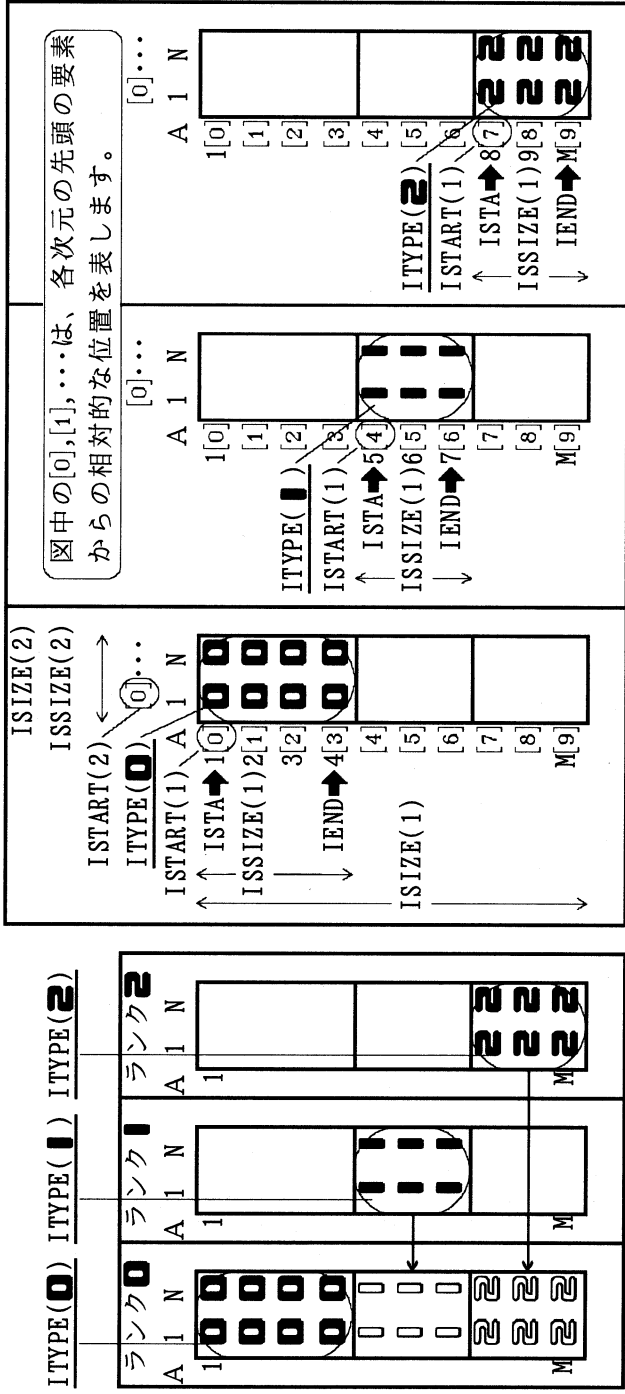


図4-6-9(1)

図4-6-9(2)

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER (M=~ , N=~ , NCPU=3)
  REAL A(M,N)
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  INTEGER ISIZE(2), ISIZE(2), ISTART(2) [1]
  INTEGER ITYPE(0:NCPU-1) [1]
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~, NPROCS, IERR)
  CALL MPI_COMM_RANK(~, MYRANK, IERR)
  ISIZE(1) = M
  ISIZE(2) = N
  ISIZE(2) = N
  ISTART(2) = 0
  DO IRANK=0, NPROCS-1
    CALL PARA_RANGE
    & (1, M, NPROCS, IRANK, I, ISTART, IEND)
    ISIZE(1) = IEND- ISTART+1
    ISTART(1) = ISTART-1
  
```

図4-6-9(3)

```

& CALL MPI_TYPE_CREATE_SUBARRAY(2, I, ISIZE,
& ISTART, ISTART, MPI_ORDER_FORTRAN,
& MPI_REAL, ITYPE(IRANK), IERR)
CALL MPI_TYPE_COMMIT(ITYPE(IRANK), IERR)[1]
ENDDO
CALL PARA_RANGE
& (1, M, NPROCS, MYRANK, I, ISTART, IEND)
:
IF (MYRANK==0) THEN
DO IRANK=1, NPROCS-1
CALL MPI_IRECV(A, 1, ITYPE(IRANK), IRANK,
1, MPI_COMM_WORLD, IREQ, IERR) [2]
CALL MPI_WAIT(IREQ, ISTATUS, IERR)
ENDDO
ELSE
CALL MPI_ISEND(A, 1, ITYPE(MYRANK), 0,
1, MPI_COMM_WORLD, IREQ, IERR) [3]
CALL MPI_WAIT(IREQ, ISTATUS, IERR)
ENDIF
:

```

## 4-6-3-2 MPI\_ALLGATHERVの代替

4-3節の『計算パターン1』で並列化した場合、各プロセスは、自分が計算した部分のデータを、他の全プロセスに送信するという通信パターンが現れます。またこのパターンは、並列化の途中でデバッグを行う場合にも使用します(4-8-1節参照)。通信の方法は前節とほぼ同様です。

### ■ (ケース4) MPI-2のMPI\_IN\_PLACEが使用でき、データが連続な場合

配列を縮小していない場合、図4-6-10(1)のように、各プロセスの配列Aの空いている部分に送信できればよいのですが、集団通信サブルーチンには『送信バッファと受信バッファがメモリー上で重なってはいけない』(3-3-5節参照)という制限があるため、図4-6-10(2)のように、新たに受信バッファ用の配列Bを設ける必要があります。しかし前述のケース1と同様に、MPI-2のMPI\_IN\_PLACE(3-8-1節参照)という変数を使用すると、この制限を回避することができます(MPI-2が使用できない環境では使用できません)。

図4-6-10(3)に示すM×Nの2次元配列Aを2次元方向に分割する場合は想定します。一般にプロセスによってデータ量が異なるため、MPI\_ALLGATHERではなくMPI\_ALLGATHERVを使用して、各プロセスのデータを他の全プロセスに送信します。図4-6-10(4)の[1]で以下の2つの配列を作成し、[2]でMPI\_ALLGATHERV(付録の「MPI\_IN\_PLACEの使い方」参照)をコールします。[2]では送信バッファの代わりにMPI\_IN\_PLACEを指定します。

● IRCNT(i) : ランクiのプロセスが担当する要素数を設定します。

● IDISP(i) : ランクiのプロセスが担当する最初の要素の、配列の先頭からの変位(⑩,⑪,⑫)を設定します。

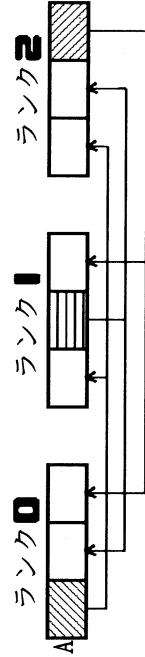


図4-6-10(1)

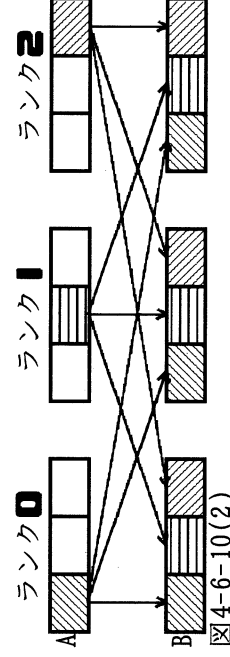


図4-6-10(2)

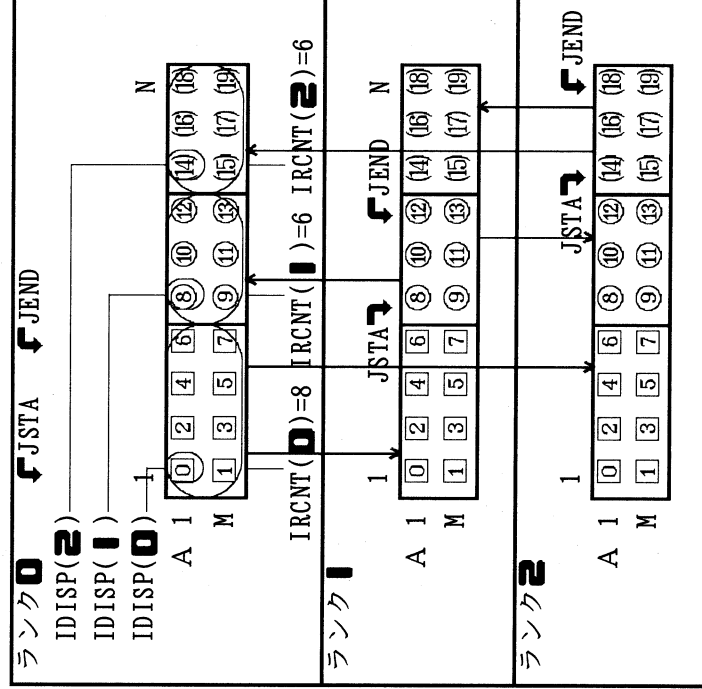


図4-6-10(3)

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER (M=~, N=~, NCPU=3)
  REAL A(M,N)

  INTEGER IRCNT(0:NCPU-1), IDISP(0:NCPU-1) [1]
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~, NPROCS, IERR)
  CALL MPI_COMM_RANK(~, MYRANK, IERR)
  DO IRANK=0, NPROCS-1
    CALL PARA_RANGE
    & (1, N, NPROCS, IRANK, JSTA, JEND)
    IRCNT(IRANK) = M*(JEND-JSTA+1)
    IDISP(IRANK) = M*(JSTA-1)
  ENDDO
  CALL PARA_RANGE
  & (1, N, NPROCS, MYRANK, JSTA, JEND)
  :
  CALL MPI_ALLGATHERV
  & (MPI_IN_PLACE, IDUMMY, IDUMMY, [2]
  & A, IRCNT, IDISP, MPI_REAL,
  & MPI_COMM_WORLD, IERR)
  :
  
```

図4-6-10(4)



■ (ケース5) MPI-2のMPI\_IN\_PLACEが使用できず、データが連続な場合

MPI-2のMPI\_IN\_PLACEの機能が使用できないマシン環境で、前述のケース4のように、配列Bを新たに設けずに配列Aの空いている部分に直接受信したい場合、MPI\_ALLGATHERVと同等の通信を、集団通信サブルーチンMPI\_BCASTを使用して行います。

まず図4-6-11(2)の[1]で以下の配列を作成します。

- JJSTA(i) : ランクiのプロセスが担当する部分の、2次元方向の最初の要素番号を設定します。
- JJLEN(i) : ランクiのプロセスが担当する要素数を設定します。

[2]がMPI\_ALLGATHERVの代替で通信を行う部分です。図4-6-11(1)から分かるように、各プロセスは自分が担当するデータを他の全プロセスに送信すればよいので、MPI\_BCASTを使用して通信を行います。[2]の下線部に示すように、送信元のプロセスが**0,1,2**と順に変化します。

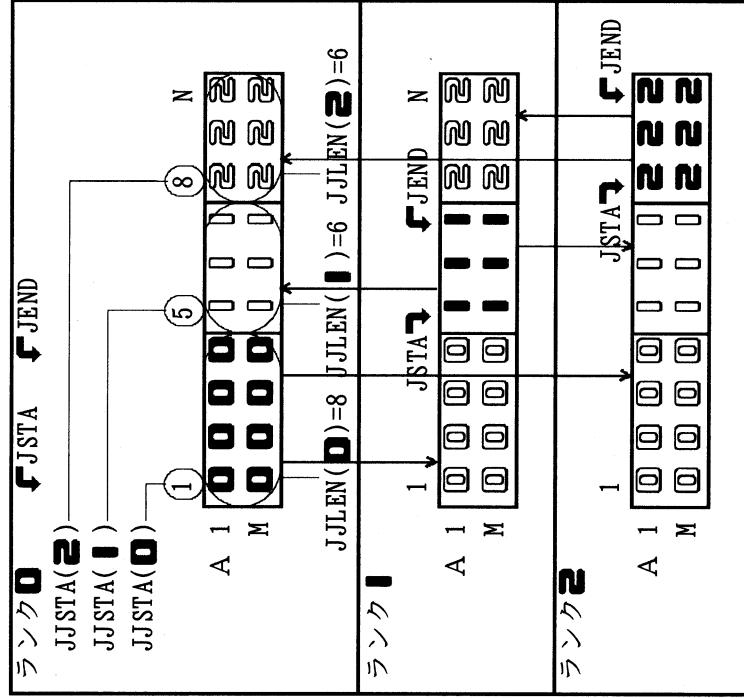


図4-6-11(1)

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER (M=~, N=~, NCPU=3)
  REAL A(M,N)
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  INTEGER JJSTA(0:NCPU-1), JJLEN(0:NCPU-1) [1]
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~, NPROCS, IERR)
  CALL MPI_COMM_RANK(~, MYRANK, IERR)
  DO IRANK=0, NPROCS-1
    CALL PARA_RANGE
    & (1, N, NPROCS, IRANK, JSTA, JEND)
    JJSTA(IRANK) = JSTA
    JJLEN(IRANK) = M*(JEND-JSTA+1) [1]
  ENDDO
  CALL PARA_RANGE
  & (1, N, NPROCS, MYRANK, JSTA, JEND)
  :
  DO IRANK=0, NPROCS-1
    CALL MPI_BCAST(A(1, JJSTA(IRANK)),
    & JJLEN(IRANK), MPI_REAL, IRANK, [2]
    & MPI_COMM_WORLD, IERR)
  ENDDO
  :
  
```

図4-6-11(2)

■ (ケース6) データが不連続で派生データ型を使う場合

ケース4とケース5では、各プロセスの送信データがメモリー上で連続していました。図4-6-12(1)では配列Aを1次元方向に分割しているので、各プロセスの送信データはメモリー上でとびとびになっています(Fortranの場合)。とびとびのデータを、いったん連続した一時配列に入れてから送信(受信はその逆)する方法もありますが、それが面倒な場合、前述のケース3と同様に、3-5節で説明した派生データ型を使用します。

図4-6-12(3)の[1]で、MPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照)(MPI-2が使用できる場合のみ)を使用して、ランクiのプロセスの担当部分の派生データ型を配列ITYPE(i)に作成します(図4-6-12(2)参照)。ケース4で使用したMPI\_ALLGATHERVでは、受信バッファのデータ型を1種類しか指定できない(つまりプロセスごとに配列で別々に指定できない)ため、ケース5と同じ方法で通信を行います。

[2]のMPI\_BCASTは、送信元のプロセス(ランクIRANK)では、『配列A(の先頭)から始まる1個のITYPE(IRANK)型のデータを、他の全てのプロセスに送信する』という意味になり、その他のプロセスでは、『ランクIRANKのプロセスから送られたデータを、配列A(の先頭)から始まる1個のITYPE(IRANK)型の受信バッファに受信する』という意味になります。

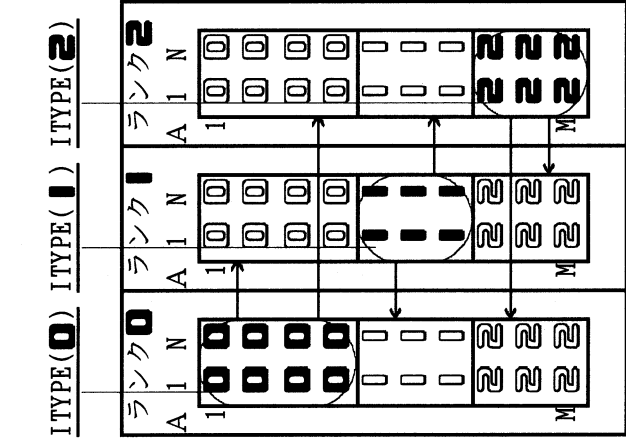


図4-6-12(1)

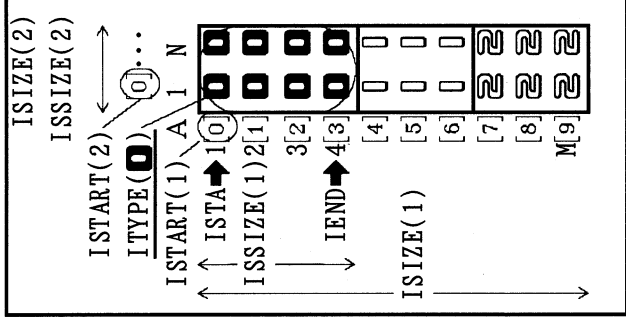
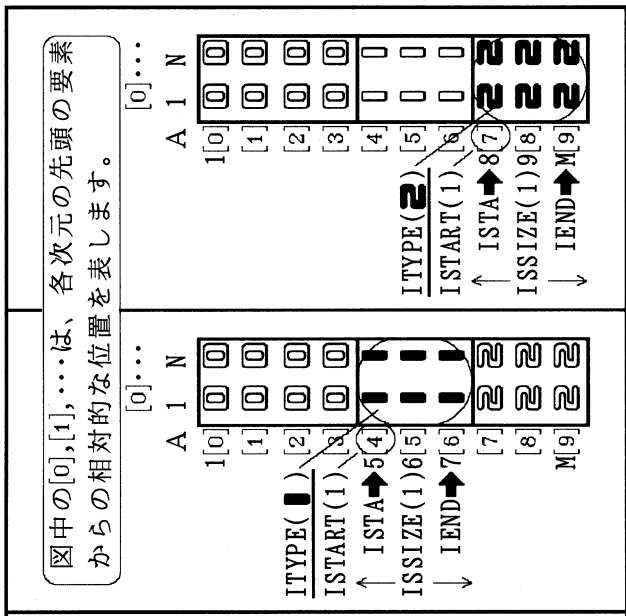


図4-6-12(2)



```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER (M=~ , N=~ , NCPU=3)
  REAL A(M,N)
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  INTEGER ISIZE(2), ISSIZE(2), ISTART(2) [1]
  INTEGER ITYPE(0:NCPU-1) [1]
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~, NPROCS, IERR)
  CALL MPI_COMM_RANK(~, MYRANK, IERR)
  ISIZE(1) = M
  ISIZE(2) = N
  ISSIZE(2) = N
  ISTART(2) = 0
  
```

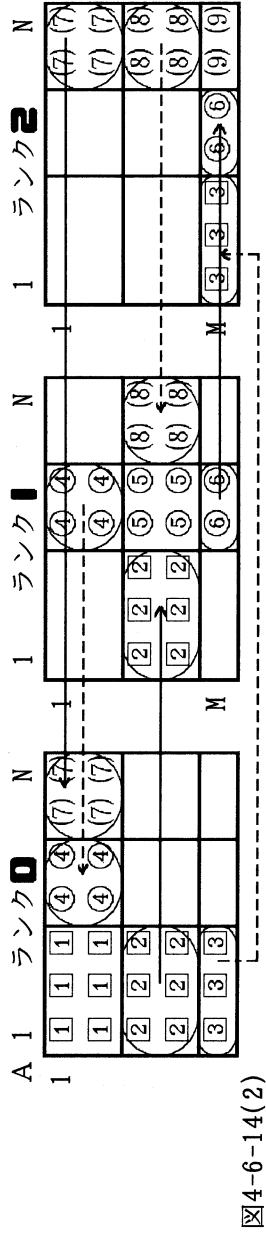
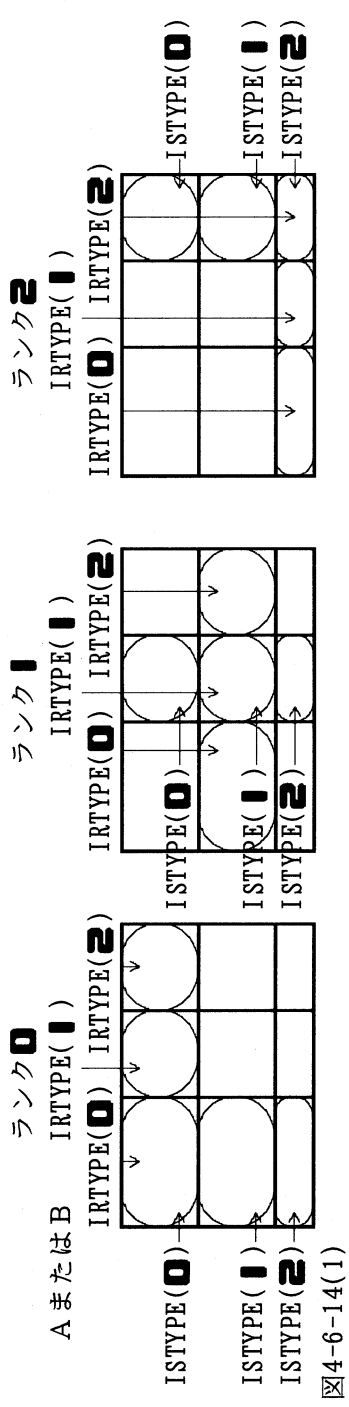
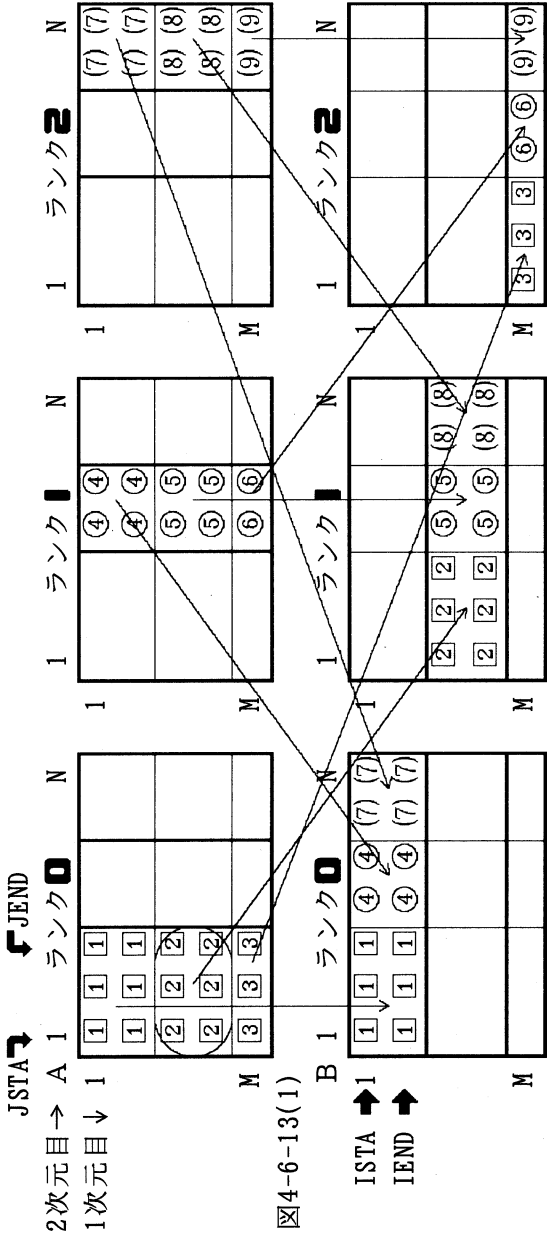
```

DO IRANK=0, NPROCS-1
  CALL PARA_RANGE
  & (1, M, NPROCS, IRANK, ISTART, IEND)
  ISSIZE(1) = IEND-ISTART+1
  ISTART(1) = ISTART-1
  CALL MPI_TYPE_CREATE_SUBARRAY(2, ISIZE,
  & ISSIZE, ISTART, MPI_ORDER_FORTRAN, [1]
  & MPI_REAL, ITYPE(IRANK), IERR)
  CALL MPI_TYPE_COMMIT(ITYPE(IRANK), IERR)
ENDDO
CALL PARA_RANGE
& (1, M, NPROCS, MYRANK, ISTART, IEND)
:
DO IRANK=0, NPROCS-1
  CALL MPI_BCAST(A, 1, ITYPE(IRANK), IRANK,
  & MPI_COMM_WORLD, IERR) [2]
ENDDO
:
  
```

図4-6-12(3)

本節で説明する通信パターンは、多次元FFT、ADI法などを並列化する場合に現れます。なお、前者は並列版の数値計算ライブラリー(6章参照)で多次元FFTのサブルーチンが提供されることが多く、後者は本節の方法よりもパイプライン法(4-6-7節)やツイスト分割法(4-6-8節)を使用した方が通信量が少なくなります。

例えば並列プログラムの前半では図4-6-13(1)のように配列Aを2次元目でブロック分割して処理し、プログラムの後半では図4-6-13(2)のように配列B(またはA)を1次元目でブロック分割して処理する場合、後半の処理を行う前に、配列AからB(またはA)へ矢印に示す通信が必要になります。



## ■ 通信方法1

図4-6-13の通信パターンは、MPIの集団通信サブルーチンMPI\_ALLTOALLVに相当します。図4-6-13では配列を縮小(4-5-7節参照)していないので、例えば○で囲んだ②の要素はメモリー上で不連続になっています。不連続なデータをいったん送(受)信バッファーに圧縮するのは面倒なので、ここでは図4-6-14(1)のような派生データ型(3-5節参照)を作成して通信で使用することにします。

図4-6-14(1)で、ISTYPE(i)はランクiのプロセスに送信するデータが入っている部分の派生データ型、IRTYPE(i)はランクiのプロセスから受信したデータが入る部分の派生データ型を示します。図から分かるように、相手プロセスごとに別々の派生データ型を用意する必要があります。ところがサブルーチンMPI\_ALLTOALLVのデータ型を指定する引数(sendtypeとrecvtype)には1種類しか指定することができません。このような場合、MPI-2で追加された集団通信サブルーチンMPI\_ALLTOALLW(付録参照)を使用します(MPI-2を使用できないマシン環境では使用できません)。

なお図4-6-13では、図4-6-14(2)のように配列Aから配列Aに通信せずに、配列Bに通信しています。この理由ですが、集団通信サブルーチンには「送信バッファーと受信バッファーは重なってはならない」という制限があります(3-3-5節参照)。配列Aから配列Aに通信した場合、図4-6-13で例えばランク②の④の部分に、送信バッファーと受信バッファーの両方に属するため、この制限に引っ掛かります。このため、配列Aとは別の配列Bを受信バッファーとして使用します(MPI\_ALLTOALLWでは、MPI-2で追加されたMPI\_IN\_PLACE(3-8節参照)も使用できません)。プログラム例を図4-6-15(1)に示します。

● 図4-6-15(1)の[1]で、図4-6-14(1)に示す派生データ型を、3-5-3-2節で説明した自作のサブルーチンPARA\_TYPE\_BLOCK2を使用して作成しMPI-2が使用できるマシン環境の場合は、MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照)を使用して下さい)、配列ISTYPEとIRTYPEに入れます。

● [2]で、MPI\_ALLTOALLWの引数(詳細は付録参照)で指定する以下の配列を設定します。

- 配列ISCNT(i)にはランクiに送信するデータ数を設定し、配列IRCNT(i)にはランクiから受信するデータ数を設定します。本例では派生データ型を使用するので、どちらも「1(個)」を設定します。

- 配列ISDISP(i)に、送信バッファーAの先頭アドレスから、ランクiに送信するデータが入っている場所までの変位(単位はバイト)を設定します。同様にIRDISP(i)に、受信バッファーBの先頭アドレスから、ランクiから受信したデータを入れる場所までの変位(単位はバイト)を設定します。派生データ型ISTYPE,IRTYPEの起点は配列A,Bの先頭アドレスなので(3-5-3節参照)、先頭アドレスを示す「0」を設定します。

● [3]を実行すると図4-6-13の矢印の通信が行われます。

## ■ 通信方法2

図4-6-13では受信バッファーとして配列Bを使用しましたが、図4-6-14(2)のように配列Bを使用せずに配列Aから配列Aに直接通信する方法を、図4-6-15(2)のプログラムで説明します。

● [5]で図4-6-15(1)と同様に図4-6-14(1)に示す派生データ型ISTYPE,IRTYPEを作成します。

● [6]のループが反復し、[7]で自分以外の全プロセスに対してデータを送信し、[8]で自分以外の全プロセスからデータを受信します。このときIREQ1,IREQ2は、[4]に示すように各プロセスごとの配列となります。

● [9]で、[7]と[8]の通信に対するMPI\_WAITを実行します。これで図4-6-14(2)の矢印に示す通信が行われます。

## ■ 通信方法3,4

上記の「通信方法2」では自分以外の全プロセスに対して一気に送受信を行いました。この方法では1つずつ順に通信を行うので、通信する相手のプロセス数が多い時は「通信方法2」よりも無難だと思われま

● 図4-6-15(3)の[10]でI=1のとき、図4-6-14(2)の太い矢印に示すように、各プロセスは自分よりランクが1つ大きいプロセスに対して送信し(ランク②はランク①に送信)、ランクが1つ小さいプロセスから受信します。

● 次に[10]でI=2となり、図4-6-14(2)の点線の矢印に示すように、ランクが2つ大きいプロセスに対して送信し、ランクが2つ小さいプロセスから受信します。

● このように[10]のIを1つずつ増やしながら、上記の送受信を自分以外の全ての相手と行います。

● 図4-6-15(3)と同様の通信をMPI\_SENDRCV(4-6-2節,付録参照)を用いて行ったのが図4-6-15(4)です。

```

:
INCLUDE 'mpif.h'
PARAMETER (M=5, N=7, NCPU=3)
REAL A(M, N), B(M, N)
INTEGER ISCNT (0:NCPU-1), IRCNT (0:NCPU-1)
INTEGER ISDISP(0:NCPU-1), IRDISP(0:NCPU-1)
① INTEGER ISTYPE(0:NCPU-1), IRTYPE(0:NCPU-1)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
CALL PARA_RANGE
& (1, M, NPROCS, MYRANK, ISTA, IEND)
CALL PARA_RANGE
& (1, N, NPROCS, MYRANK, JSTA, JEND)
DO IRANK=0, NPROCS-1
CALL PARA_RANGE(1, M, NPROCS, IRANK,
IISTA, I IEND)
CALL PARA_TYPE_BLOCK2(1, M, 1, I ISTA,
I IEND, JSTA, JEND, MPI_REAL, I ITYPE)
& ISTYPE(IRANK) = I ITYPE
CALL PARA_RANGE(1, N, NPROCS, IRANK,
J JSTA, J JEND)
CALL PARA_TYPE_BLOCK2(1, M, 1, I ISTA, I IEND,
J JSTA, J JEND, MPI_REAL, I ITYPE)
& IRTYPE(IRANK) = I ITYPE
ENDDO
DO IRANK=0, NPROCS-1
ISCNT (IRANK) = 1
IRCNT (IRANK) = 1
ISDISP(IRANK) = 0
IRDISP(IRANK) = 0
ENDDO
:
CALL MPI_ALLTOALLW(A, ISCNT, ISDISP, I ISTYPE,
& B, IRCNT, IRDISP, I IRTYPE, [3]
& MPI_COMM_WORLD, I IERR)
:

```

図4-6-15(1) 通信方法1

```

:
DO I=1, NPROCS-1
IUP = MOD(MYRANK +I, NPROCS)
IDOWN = MOD(MYRANK+NPROCS-I, NPROCS)
CALL MPI_SENDECV
& (A, 1, ISTYPE(IUP) , IUP , 1,
& A, 1, IRTYPE(IDOWN), IDOWN, 1,
& MPI_COMM_WORLD, ISTATUS, IERR)
ENDDO
:

```

図4-6-15(4) 通信方法4

```

:
INCLUDE 'mpif.h'
PARAMETER (M=5, N=7, NCPU=3)
REAL A(M, N)
IREQ1(0:NCPU-1), IREQ2(0:NCPU-1)
INTEGER ISTATUS(MPI_STATUS_SIZE)
左図の①と同じ
:
DO IRANK=0, NPROCS-1
IF (IRANK/=MYRANK) THEN
CALL MPI_ISEND(A, 1, ISTYPE(IRANK),
↑ IRANK, 1, MPI_COMM_WORLD,
[7] IREQ1(IRANK), IERR)
* CALL MPI_IRECV(A, 1, IRTYPE(IRANK),
[8] IRANK, 1, MPI_COMM_WORLD,
IREQ2(IRANK), IERR)
↓ ENDDO
DO IRANK=0, NPROCS-1
IF (IRANK/=MYRANK) THEN
CALL MPI_WAIT
& (IREQ1(IRANK), ISTATUS, IERR)
[9] CALL MPI_WAIT
& (IREQ2(IRANK), ISTATUS, IERR)
ENDDO
:

```

図4-6-15(2) 通信方法2

```

:
INCLUDE 'mpif.h'
PARAMETER (M=5, N=7, NCPU=3)
REAL A(M, N)
INTEGER ISTATUS(MPI_STATUS_SIZE)
左図の①と同じ
:
DO I=1, NPROCS-1
IUP = MOD(MYRANK+I , NPROCS)
IDOWN = MOD(MYRANK-I+NPROCS, NPROCS)
CALL MPI_ISEND(A, 1, ISTYPE(IUP),
& IUP , 1, MPI_COMM_WORLD, IREQ1, IERR)
CALL MPI_IRECV(A, 1, IRTYPE(IDOWN),
& IDOWN, 1, MPI_COMM_WORLD, IREQ2, IERR)
CALL MPI_WAIT(IREQ1, ISTATUS, IERR)
CALL MPI_WAIT(IREQ2, ISTATUS, IERR)
ENDDO
:

```

図4-6-15(3) 通信方法3

## 4-6-4 合計(内積)/最大(最小)

図4-6-16(1)を並列化したプログラムを図4-6-16(2)に示します。

図4-6-16(2)の最初のループでは、②で合計を、③で内積を求めていきます。このように、並列化したループの中に合計や内積を求める計算が含まれている場合、ループを終了した直後に、「各プロセスが求めた合計の総合計」を求めるために、MPI\_SUMを使用して、MPI\_REDUCEまたはMPI\_ALLREDUCEで通信を行う必要があります。②と③は、⑤で全プロセスがSUM1とSUM2を参照するので、MPI\_ALLREDUCEで通信を行い、結果を全プロセスに送信します。またSUM1とSUM2を別々にMPI\_ALLREDUCEで通信すると、4-1節で述べたように通信の立上り時間によるオーバーヘッドが2回かかってしまうので、⑦の作業配列WORKSとWORKRを使用し、⑨では1回で通信します。

図4-6-16(2)の④で最大値を求めていきます。このように、並列化したループの中に、最大(最小)値を求める計算が含まれている場合、ループを終了した直後に、「各プロセスが求めた最大(最小)値の全体の最大(最小)値」を求めるために、MPI\_MAX(MPI\_MIN)を使用して、MPI\_REDUCEまたはMPI\_ALLREDUCEで通信を行う必要があります。④は⑥で(例えば)ランク0のプロセスだけが結果を得ればよいので、MPI\_ALLREDUCEではなく、⑩のようにMPI\_REDUCEで通信します。

4章では合計を求めるプログラムを並列化の例題として取り上げました。しかし実際のプログラムでは、合計、内積、最大値(最小値)を求めるだけのループが独立しているというのはいずれも、①のように派手な計算の片隅(通常はループの最後に)、②,③,④が『さりげなく』現れているのが一般的です。しかしループを⑧のように並列化したとき、そのループの中に②,③,④が含まれているのを見落とすと、結果がおかしくなりますので注意して下さい。

```

:
SUM1 = 0.0
SUM2 = 0.0
AMAX = 0.0
DO I = 1, N
  A(I) = A(I) + B(I)
  C(I) = C(I) + D(I)
  E(I) = E(I) + F(I)
  G(I) = G(I) + H(I)
  X(I) = X(I) + Y(I)
  SUM1 = SUM1 + A(I)
  SUM2 = SUM2 + P(I)*Q(I)
  IF (A(I) > AMAX) AMAX = A(I)
ENDDO
DO I = 1, N
  G(I) = G(I)*SUM1 + SUM2
ENDDO
PRINT *,AMAX
:

```

図4-6-16(1)

```

:
INCLUDE 'mpif.h'
DIMENSION WORKS(2),WORKR(2)
SUM1 = 0.0
SUM2 = 0.0
AMAX = 0.0
DO I = ISTART, IEND
  A(I) = A(I) + B(I)
  C(I) = C(I) + D(I)
  E(I) = E(I) + F(I)
  G(I) = G(I) + H(I)
  X(I) = X(I) + Y(I)
  SUM1 = SUM1 + A(I)
  SUM2 = SUM2 + P(I)*Q(I)
  IF (A(I) > AMAX) AMAX = A(I)
ENDDO
WORKS(1) = SUM1
WORKS(2) = SUM2
CALL MPI_ALLREDUCE(WORKS,WORKR,2,MPI_REAL,
& MPI_SUM,MPI_COMM_WORLD,IERR)
SUM1 = WORKR(1)
SUM2 = WORKR(2)
CALL MPI_REDUCE(AMAX,AMAX,1,MPI_REAL,
& MPI_MAX,0,MPI_COMM_WORLD,IERR)
AMAX = AMAX
DO I = ISTART, IEND
  G(I) = G(I)*SUM1 + SUM2
ENDDO
IF (MYRANK==0) PRINT *,AMAX
:

```

図4-6-16(2)

一方、並列化するループの外にタイムステップのループがある場合、合計(または内積)を求める演算には図4-6-16(3)のような3つのパターンがあります。

- ①ではタイムステップごとに内側のループの合計を求め、結果を書き出しています。
- ②ではタイムステップ全体の合計を求め、結果を最後に1回だけ書き出しています。
- ③ではタイムステップごとに、タイムステップの最初からその時点までの累計を求め、結果を書き出しています。

内側のループを並列化した場合、それぞれのパターンに応じて図4-6-16(4)のように修正します。

- ①では内側のループが終了するたびにMPI\_REDUCEを行います。
- ②ではタイムステップループが終了した後1回だけMPI\_REDUCEを行います。
- ③では少し注意が必要です。まず内側のループに入る前にSUM3に入っているその時点までの累計値を一時変数TEMPに保管します。そしてSUM3をゼロクリアしてから内側のループを実行します。ループ終了後にMPI\_REDUCEを行い、その結果が入るSSUM3と、TEMPに保管した値を合計してSUM3に戻します。

```

:
SUM2 = 0.0           ②
SUM3 = 0.0           ③
D0 タイムステップ・ループ
:
SUM1 = 0.0           ①
D0 I = 1, N
SUM1 = SUM1 + A(I) ①
SUM2 = SUM2 + B(I) ②
SUM3 = SUM3 + C(I) ③
ENDDO
PRINT *,SUM1,SUM3 ①③
:
PRINT *,SUM2       ②
:

```

図4-6-16(3)



```

:
INCLUDE 'mpif.h'
:
SUM2 = 0.0           ②
SUM3 = 0.0           ③
D0 タイムステップ・ループ
:
SUM1 = 0.0           ①
TEMP = SUM3          ③
SUM3 = 0.0           ③
D0 I = ISTART, IEND
SUM1 = SUM1 + A(I) ①
SUM2 = SUM2 + B(I) ②
SUM3 = SUM3 + C(I) ③
ENDDO
CALL MPI_REDUCE(SUM1,SSUM1,~,) ①
SUM1 = SSUM1         ①
CALL MPI_REDUCE(SUM3,SSUM3,~,) ③
SUM3 = TEMP + SSUM3 ③
IF (MYRANK==0) PRINT *,SUM1,SUM3 ①③
:
ENDDO
CALL MPI_REDUCE(SUM2,SSUM2,~,) ②
SUM2 = SSUM2         ②
IF (MYRANK==0) PRINT *,SUM2 ②
:

```

図4-6-16(4)

## 4-6-5 特定の配列要素の参照

本例は、4-6-1節の「メッセージ交換はどんなときに必要か」の例と同じです。実際のプログラムではこのような例はめったに現れませんが、もし現れた場合には以下のような対応をとる必要がありますので注意して下さい。

図4-6-17(1)の下線部に、配列の添字にループ反復『I』以外の値(この例では『1』)が現れています。まず配列を縮小せずに並列化した図4-6-17(2)では、①のループを並列化したため、ループを抜けた直後にはランク0のプロセスのみがA(1)の値を持っています。③で全プロセスがA(1)の値を参照しますが、ランク1と2のプロセスはA(1)の値を持っていないため、③を行う前に②でランク0からその他のプロセスにA(1)の値を送信します(図4-6-17(4)参照)。

一方、配列を縮小して並列化した図4-6-17(3)では、ランク1と2のプロセスはA(1)に自分の担当する要素の値を持っています。このため、②のようにランク0からその他のプロセスにA(1)の値を送信すると、図4-6-17(5)に示すように、ランク1と2のプロセスでは自分の担当するA(1)の値が壊されてしまいます(図中の×)。これを防ぐため、ランク0のプロセスはA(1)の値を④で一時変数TEMPに入れ、それを⑤でその他のプロセスに送信し、⑥ではTEMPの値を参照するようになります。

```

:
:
DIMENSION A(9),B(9)
:
:
DO I = 1, 9
  A(I) = ~
ENDDO
DO I = 1, 9
  B(I) = B(I)*A(1)
ENDDO
:

```

図4-6-17(1)

```

:
:
DIMENSION A(9),B(9)
:
:
DO I = ISTART, IEND
  A(I) = ~
ENDDO
CALL MPI_BCAST(A(1),1,MPI_REAL,
& 0,MPI_COMM_WORLD,IERR) ...
DO I = ISTART, IEND
  B(I) = B(I)*A(1)
ENDDO
:

```

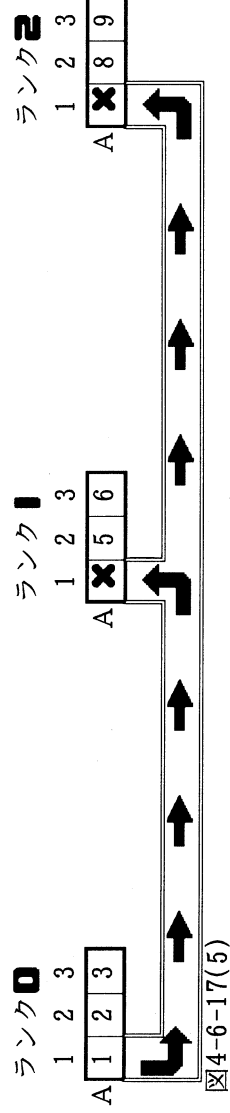
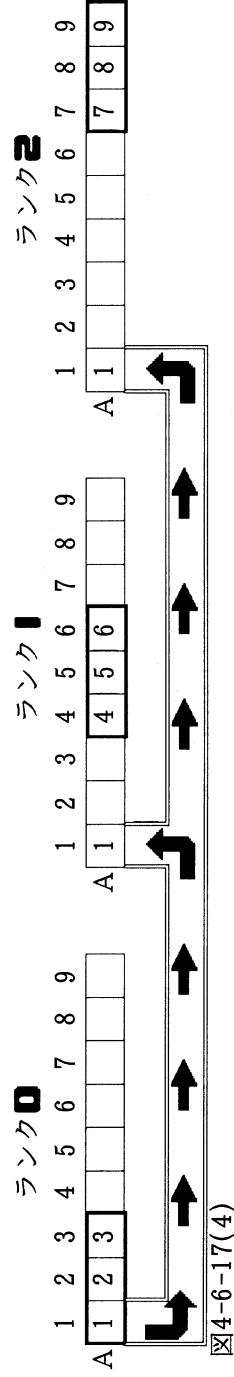
図4-6-17(2)

```

:
:
DIMENSION A(3),B(3)
:
:
DO I = 1, IEND-ISTART+1
  A(I) = ~
ENDDO
IF (MYRANK == 0) TEMP = A(1)
CALL MPI_BCAST(TEMP,1,MPI_REAL,
& 0,MPI_COMM_WORLD,IERR)
DO I = 1, IEND-ISTART+1
  B(I) = B(I)*TEMP
ENDDO
:

```

図4-6-17(3)





通常、MPI\_(ALL)REDUCEなどの通信しながら演算を行うサブルーチンは、並列化したループ内に合計、内積、最大(小)値などを求める計算が含まれている場合、ループ終了後に「合計の合計」や「最大の最大」を求めるために使用しますが(4-6-4節参照)、そのほかにも便利な使用方法があります。

### 4-6-6-1 データの収集

図4-6-18(1)のD0ループを並列化する際、何らかの理由で、1~6の各ループ反復を図4-6-18(2)の配列MAP内のランクが担当するとします。またMAP内のランク値は不規則であります(計算領域が非構造格子で構成されているプログラムなどでこのパターンが現れます)。

図4-6-18(3)にMAPを使って並列化したプログラムを示します。D0ループを3プロセスで並列に実行し、D0ループを抜けた直後、配列Xの自身は図4-6-18(4)の上部のようになります。

次に、各プロセスが計算した配列Xの要素をランク0に収集します。「正しい」(通信量が最も少ない)収集方法を図4-6-18(4)の下部に示します。まず各プロセスは自分が担当した部分のデータを圧縮し(一般にメモリー上でとびとびになっているため)、それをMPI\_GATHERVなどでランク0に収集し、ランク0が配列Xの元の位置に復元します。しかしこのようにデータを圧縮、復元するのは若干面倒です。

このような場合、以下で説明する重ね合わせを用いると、圧縮、復元をせずに簡単にデータの収集を行うことができます。ただし、後述するように通信時間がかかります。

なお、3-5-2節で紹介したMPI\_TYPE\_CREATE\_INDEXED\_BLOCKを使って、とびとびのデータの派生データ型を作成すると、圧縮、復元は不要になり、通信時間も短くて済みますが、以下では説明は省略します。

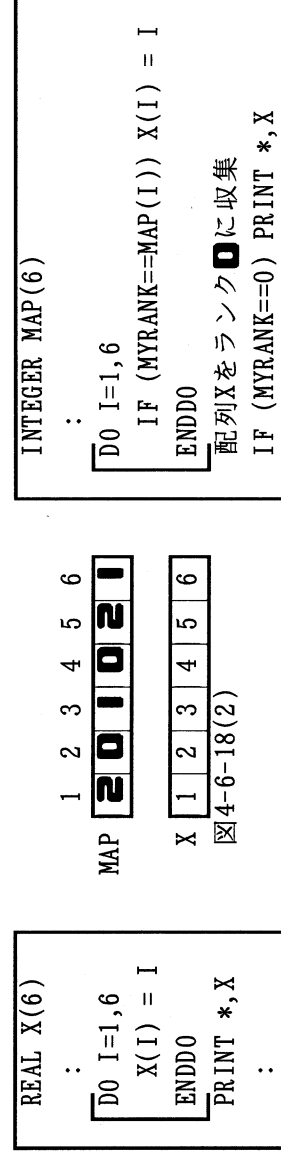


図4-6-18(1)

図4-6-18(3)

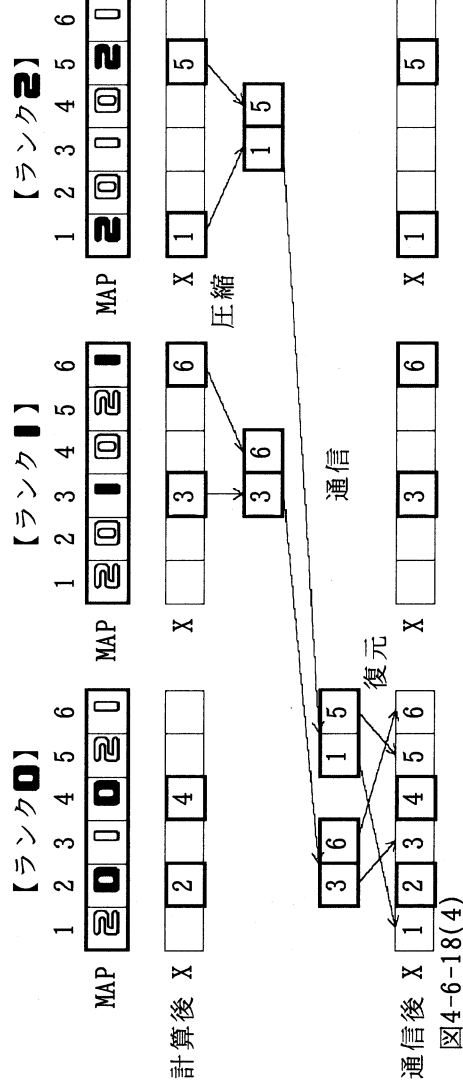


図4-6-18(4)

重ね合わせを使用したプログラムを図4-6-18(5)に、データの動きを図4-6-18(7)に示します。各プロセスは、図4-6-18(5)の①で配列X全体をゼロクリア(これが前提になります)してから、②で自分の担当部分を計算します。計算終了後、③で配列X全体をMPI\_REDUCEで通信すると、配列Xの各要素が要素ごとに加算されてランク0の配列XXに入り、図4-6-18(4)と同じ結果を簡単に得ることができます。

なお3-3-5節で説明したように、MPI\_REDUCEでは、送信バッファXと受信バッファXXは別の配列にしななければならないので注意して下さい(ただし、MPI-2で追加されたMPI\_IN\_PLACEという変数を使用すると、この制限を回避することができます。詳細は3-8-1節を参照して下さい)。

重ね合わせでは、値の入っていない部分のゼロクリアが前提になります。図4-6-18(5)では計算前に①で配列X全体をゼロクリアしましたが、図4-6-18(6)の④に示すように、通信前に配列Xの自分の担当部分以外の要素のみをゼロクリアする方法もあります。

```

:
DO I=1,6
  X(I) = 0.0
ENDDO
DO I=1,6
  IF (MYRANK==MAP(I)) X(I) = I
ENDDO
CALL MPI_REDUCE(X,XX,6,MPI_REAL,MPI_SUM,0,
& MPI_COMM_WORLD,IERR) ③
IF (MYRANK==0) PRINT *,XX
:

```

図4-6-18(5)

```

:
DO I=1,6
  IF (MYRANK==MAP(I)) X(I) = I
ENDDO
DO I=1,6
  IF (MYRANK/=MAP(I)) X(I) = 0.0
ENDDO
CALL MPI_REDUCE(X,XX,6,MPI_REAL,MPI_SUM,0,
& MPI_COMM_WORLD,IERR)
IF (MYRANK==0) PRINT *,XX
:

```

図4-6-18(6)

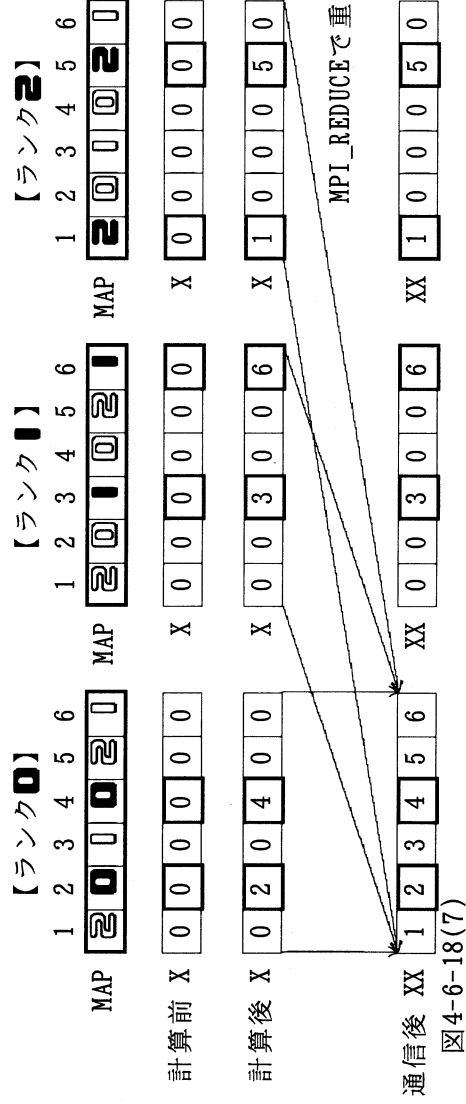
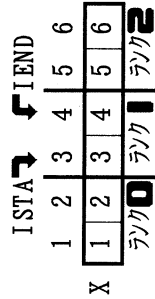


図4-6-18(7)

図4-6-19(1)～(3)のように、ブロック分割したループに対して重ね合わせを行った場合、ゼロクリアは⑤(計算前)、または⑥(通信前)のように行います。

重ね合わせは図4-6-18(4)の「正しい」収集方法よりも簡単なので、例えば並列プログラムのテスト段階で、とりあえず各プロセスが計算した結果を簡単に収集して単体プログラムの結果と比較したい場合などに用いることができます。

ただし、余分なデータ(ゼロ)を通信するため、図4-6-18(4)の「正しい」通信方法に比べて通信時間がかかります。本番プログラムの通信時間がクリティカルな部分では用いないで下さい。



(注) 上図の ISTA, IEND は  
ランク 用

図4-6-19(1)

```

:
DO I=1,6
  X(I) = 0.0
ENDDO
DO I=ISTA, IEND
  X(I) = I
ENDDO
CALL MPI_REDUCE(~)
:

```

⑤

図4-6-19(2)

```

:
DO I=ISTA, IEND
  X(I) = I
ENDDO
DO I=1, ISTA-1
  X(I) = 0.0
ENDDO
DO I=IEND+1, 6
  X(I) = 0.0
ENDDO
CALL MPI_REDUCE(~)
:

```

⑥

図4-6-19(3)

## 4-6-6-2 間接アドレスを使用した計算

重ね合せの別の例を図4-6-20(1)(2)で説明します。図4-6-20(1)の配列AとXに値が入っており、配列Aの値を、配列X(のインデックス配列INDで示す場所)に加算します。

このように間接アドレスを使用する計算は、有限要素法(陽解法)で要素の値から節点の値を求める計算(またはその逆)などで現れます。

配列AとXを図4-6-20(3)に示すようにブロック分割して並列化する場合、配列Aを配列Xに加える矢印がプロセスの境界を越えているので、変数Aの境界を越えたデータを他のプロセスに通信する必要があります。まず各プロセスは配列Aのうち他のプロセスに渡すデータをプロセス別に圧縮し(一般にメモリー上でとびとびになっているため)、それをMPI\_ALLTOALLVなどで他の各プロセスに送信し、逆に他の各プロセスから送信されたデータを受信して自分の配列Aの元の位置に復元し、最後に配列Aを配列Xに加えます。しかしこのようにデータを圧縮、復元するのは若干面倒です。このような場合、重ね合わせを用いると簡単にデータの通信を行うことができます。

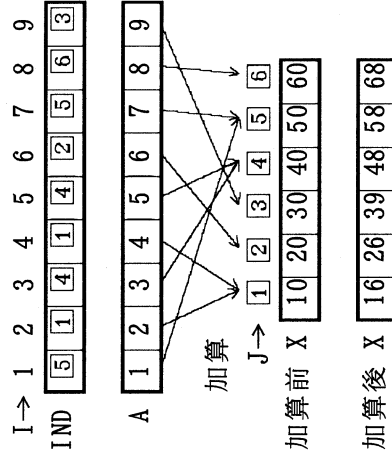


図4-6-20(1)

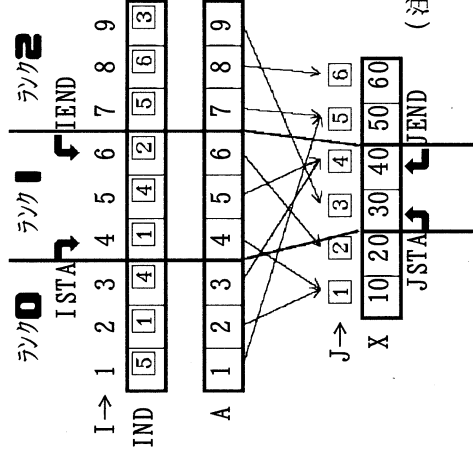


図4-6-20(2)

(注) 図の ISTA, IEND, JSTA, JEND はラック ■ 用

```

REAL A(9), X(6)
INTEGER IND(9)
:
DO I=1,9
  X(IND(I)) = X(IND(I)) + A(I)
ENDDO
PRINT *, X
:

```

図4-6-20(2)

```

INTEGER IRCNT(0:2)
:
DO J=1, JSTA-1
  X(J) = 0.0
ENDDO
DO J=JEND+1,9
  X(J) = 0.0
ENDDO
DO I=ISTA, IEND
  X(IND(I)) = X(IND(I)) + A(I)
ENDDO
DO IRANK=0,2
  IRCNT(IRANK) = 2
ENDDO
CALL MPI_REDUCE_SCATTER
& (X, XX, IRCNT, MPI_REAL, MPI_SUM,
& MPI_COMM_WORLD, IERR)
PRINT *, MYRANK, XX
:

```

図4-6-20(4)

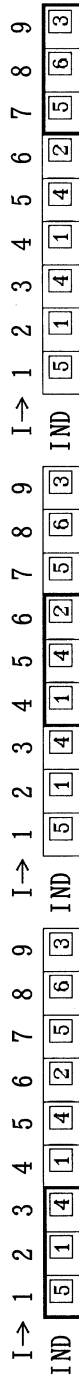
重ね合せを使用したプログラムを図4-6-20(4)に、データの動きを図4-6-20(6)に示します。各プロセスは図4-6-20(4)の①で、配列Xのうち自分の担当以外の部分をゼロクリア(これが前提になります)し、②で配列Aの自分の担当部分を配列Xに加えます。

計算終了後、④でMPI\_REDUCE\_SCATTER(MPI\_REDUCEを行った後、MPI\_SCATTERを行う集団通信ルーチン:付録参照)を使用して通信を行ないます。まずMPI\_REDUCEで配列Xの全要素が要素ごとに加算され、その結果図4-6-20(6)の点線部分のようになります。次に点線部分がMPI\_SCATTERで各プロセスに2要素ずつ(③で指定)分配されて各プロセスの配列XXに入り、図4-6-20(5)と同じ結果を簡単に得ることができます。

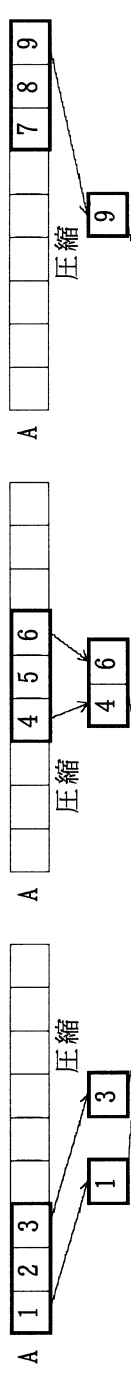
なお、3-3-5節で説明したように、MPI\_REDUCE\_SCATTERでは、送信バッファA-Xと受信バッファA-XXは別の配列にしなければならぬので注意して下さい(ただし、MPI-2で追加されたMPI\_IN\_PLACEという変数を使用すると、この制限を回避することができます。詳細は3-8-1節を参照して下さい)。

前述のように、重ね合せは「正しい」通信と比較して通信時間がかかるので、通信時間がクリティカルな部分では使用しないで下さい。

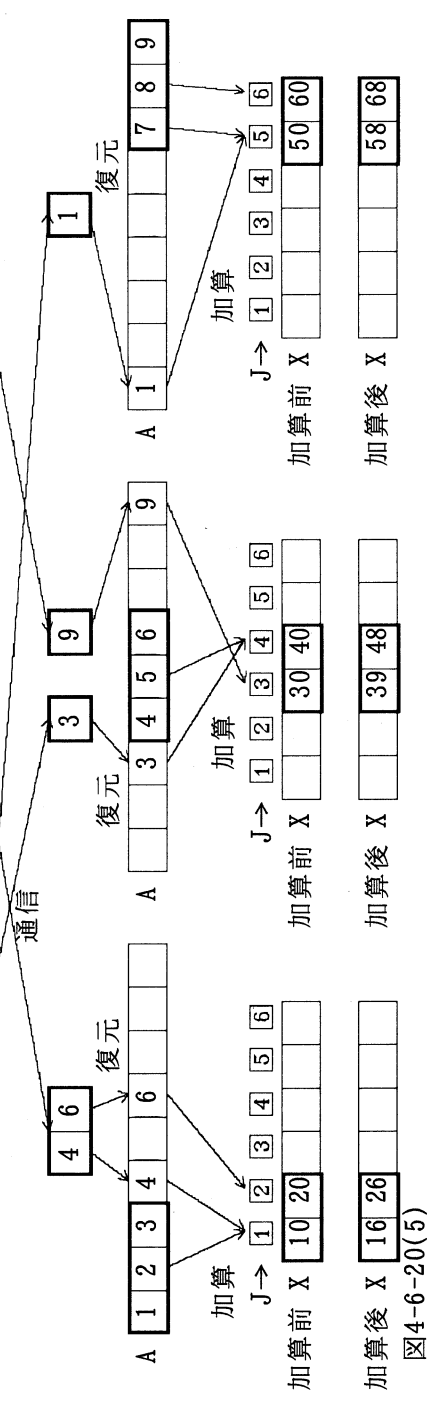
【ランク0】



【ランク1】



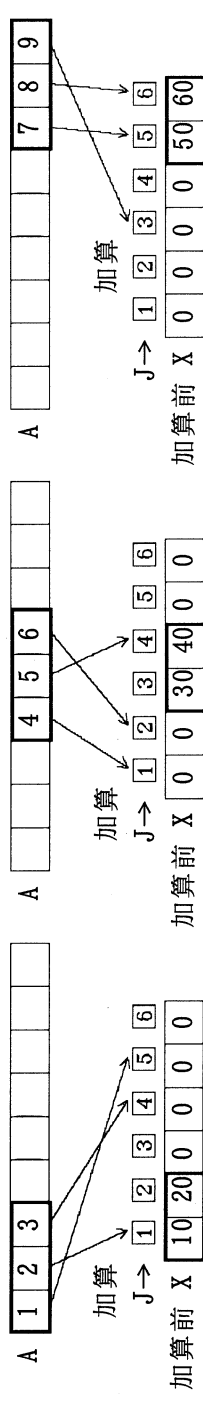
【ランク2】



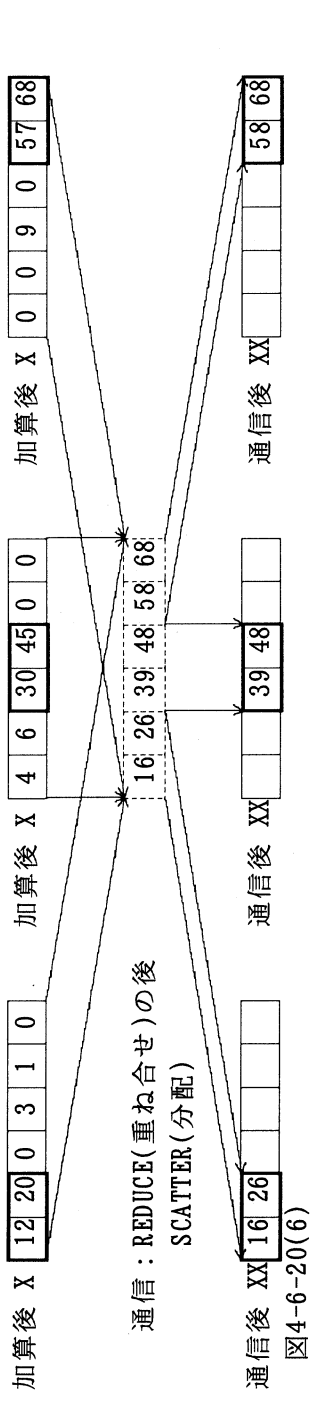
【ランク0】



【ランク1】



【ランク2】



■ パイプライン法とは

本節で説明する方法は、実際のプログラムではそれほど現れることはないのですが、お急ぎの方は飛ばしてもかまいません。

図4-6-22(1)(2)では、計算領域の各要素に対し、上の要素と左の要素から自分を計算するという依存関係があり、簡単には並列化できません(なお、『境』は境界上のデータを表します)。例えば図4-6-22(2)を2次元方向にブロック分割したとします。このときランク0と1が行う■と■の計算がすべて終了した後でないとして、ランク2は★を計算することができません。また★の計算には●が必要なので、ランク2のブロセスは★の計算の前にランク1のブロセスから●を送信してもらわなければなりません。

このパターンはパイプライン法(何とか)並列化することができます。このパターンは実際には対称帯行列の反復解法であるICCG法の前進消去、後退代入などで発生します(5-4-1節参照)。

一方、図4-6-23(1)の①のような『普通の』ループが続いていて、これを図4-6-23(2)のように2次元方向に分割して並列処理していたところへ、突然②のように横方向に依存関係があるループが現れたような場合にも、②のループに対してパイプライン法を適用することができます。

```

PROGRAM MAIN
PARAMETER (MX=7, MY=6)
DIMENSION X(0:MX, 0:MY)
:
DO J=1, MY
  DO I=1, MX
    X(I, J) = X(I, J) + X(I-1, J) + X(I, J-1)
  ENDDO
ENDDO
:
    
```

図4-6-22(1)

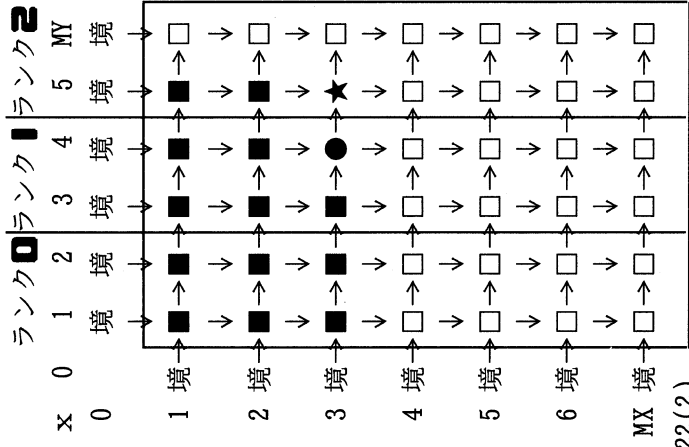


図4-6-22(2)

```

PROGRAM MAIN
PARAMETER (MX=6, MY=6)
DIMENSION X(MX, 0:MY)
:
DO J=1, MY
  DO I=1, MX
    X(I, J) = ~
  ENDDO
ENDDO
:
DO J=1, MY
  DO I=1, MX
    X(I, J) = X(I, J) + X(I, J-1)
  ENDDO
ENDDO
:
    
```

図4-6-23(1)

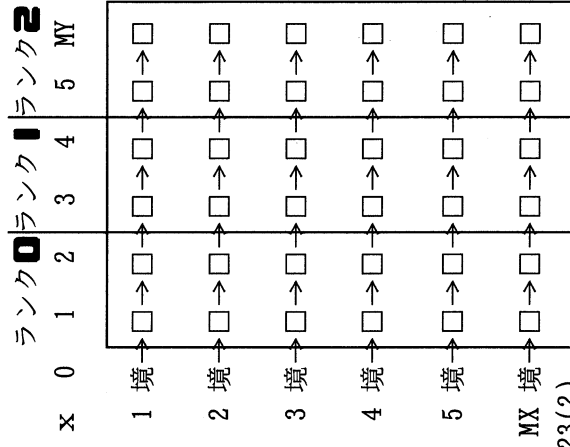


図4-6-23(2)

## ■ パイプライン法のプログラム例

図4-6-22(1)をパイプライン法で並列化した場合、配列Xを1次元方向と2次元方向のどちらでブロック分割することもできますが、ここでは図4-6-24に示すように2次元方向にブロック分割したとします。それとは別に1次元方向を適当な大きさ (IBLOCK) に分割します。分割された領域(例えば図4-6-24の着色部分)を ブロックと呼び、これが一度に処理を行う単位となります。

3つのプロセスで並列に実行した場合の動作概要を説明します。図4-6-24の( )内の数字の順に斜線が移動していき、斜線上の各ブロックが並列に処理されるというイメージになります。つまり各プロセスは隣のプロセスと1つだけ 離れたブロックを並列に処理 します。

- まず(1)に示すように、ランク**0**のプロセスのみが動作し、**1**のブロックを計算します。
- 次に(2)に示すように、ランク**0, 1**の各プロセスがそれぞれ**3**のブロックを計算します。
- 次に(3)に示すように、ランク**0, 1, 2**の各プロセスがそれぞれ**5**のブロックを計算します。この時点で初めて全プロセスが稼働します。以下同様に処理が行われます。

この処理に伴い境界の要素の通信が必要になります。例えば図4-6-24でランク**1**のプロセスは、着色したブロックを計算するためにランク**0**のプロセスが所有する**0**の要素が必要になります。従ってランク**1**のプロセスは、ランク**0**のプロセスが**1**の要素を計算した後、**0**の要素を送ってもらい、その後で**3**の計算を行う必要があります。

図4-6-22(1)を並列化したプログラムを図4-6-25に、プロセスごとのデータの動きを図4-6-26に示します。

- 図4-6-25の[1]で2次元方向にブロック分割し、各プロセスの担当範囲JSTA, JENDを決定します。
- [3]でブロックの1次元方向の大きさをIBLOCKに設定します。後述するようにIBLOCKの大きさはパフォーマンスに影響し、最適な大きさは試行錯誤で決定します。
- [4]のルーブ指標IIはこれから処理するブロックの1次元目の先頭アドレスを表し、図4-6-26の**▲**に示すように順に移動します。IIが1つ決まると、そのブロックの1次元方向の大きさ(図4-6-26の**↓**)を[5]でIBLKLENに設定します。IBLKLENは基本的にはIBLOCKと同じ大きさですが、図4-6-26のようにMXがIBLOCKで割りきれない場合、最後の半端な部分ではIBLOCKより小さくなるのでそれを調整します。
- 次に実際の計算部分について説明します。図4-6-26でランク**1**が行う着色したブロックの処理に着目すると、1つのブロックの処理は以下の3つの部分から構成されています。この処理を図4-6-25の[6]~[8]で行います。
  - (1) 左のプロセス(ランク**0**)から送信された要素を、ブロックの左の境界を越えた1列( **0** で囲んだ**②**)に受信します(図4-6-25の[6])。
  - (2) ブロック内の**3**を計算します(図4-6-25の[7])。
  - (3) ブロックの右の境界の1列( **0** で囲んだ**3**)を、右のプロセス(ランク**2**)に送信します(図4-6-25の[8])。

このように両端以外のプロセスでは、1つのブロックを受信、計算、送信の順に処理しますが、両端のプロセスでは送信もしくは受信のどちらから一方を行います(ランク**0**のプロセスでは送信のみ、ランク**2**のプロセスでは受信のみ)。これを図4-6-25の[2]で調整します。

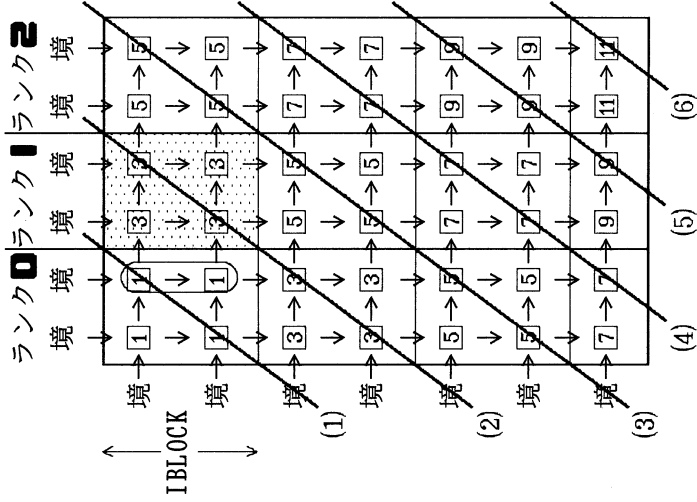


図4-6-24

```

PROGRAM MAIN
INCLUDE 'mpif.h'
PARAMETER(MX=7, MY=6)
DIMENSION X(0:MX, 0:MY)
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(ierr)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, ierr)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, ierr)
CALL PARA_RANGE(1, MY, NPROCS, MYRANK, JSTA, JEND) [1]
IUP = MYRANK+1
IF (IUP==NPROCS) IUP=MPI_PROC_NULL
IDOWN = MYRANK-1
IF (IDOWN==-1) IDOWN=MPI_PROC_NULL
IBLOCK=2
:
DO II=1, MX, IBLOCK [4]
  IBLKLEN = MIN(IBLOCK, MX-II+1) [5]
  CALL MPI_Irecv(X(II, JSTA-1), IBLKLEN, MPI_REAL, [6]
    & IDOWN, 1, MPI_COMM_WORLD, IREQR, ierr)
  CALL MPI_WAIT(IREQR, ISTATUS, ierr)
  DO J=JSTA, JEND
    DO I=II, II+IBLKLEN-1
      X(I, J) = X(I-1, J) + X(I, J-1) [7]
    ENDDO
  ENDDO
  CALL MPI_Isend(X(II, JEND), IBLKLEN, MPI_REAL, [8]
    & IUP, 1, MPI_COMM_WORLD, IREQS, ierr)
  CALL MPI_WAIT(IREQS, ISTATUS, ierr)
ENDDO
CALL MPI_FINALIZE(ierr)
:

```

図4-6-25

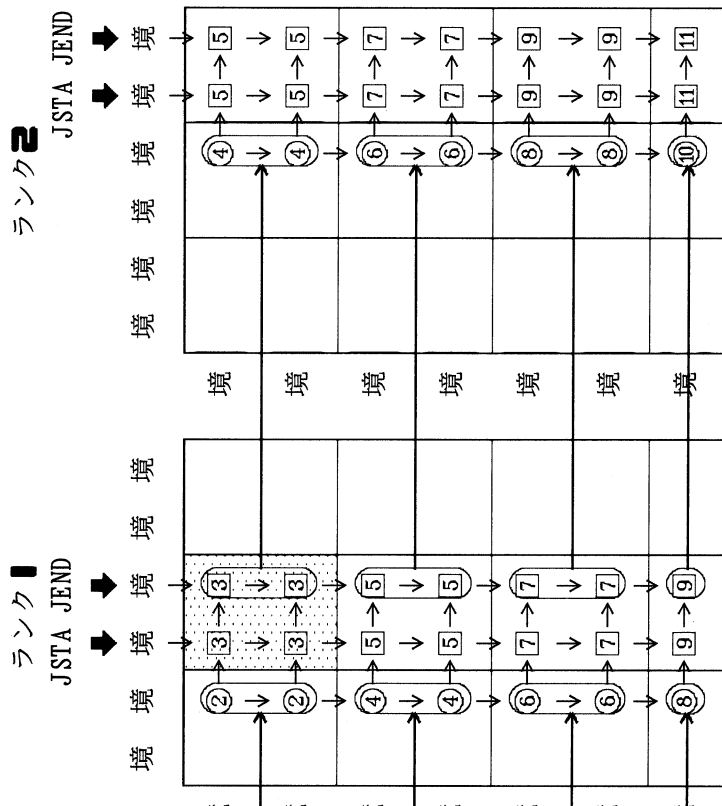
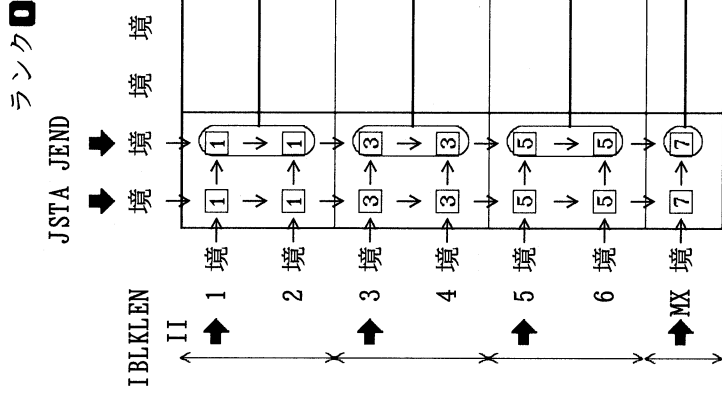
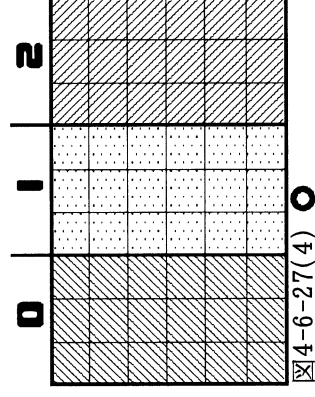
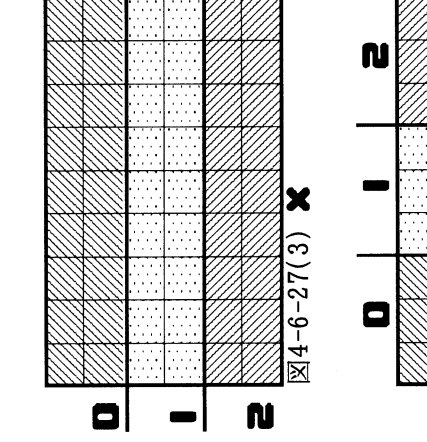
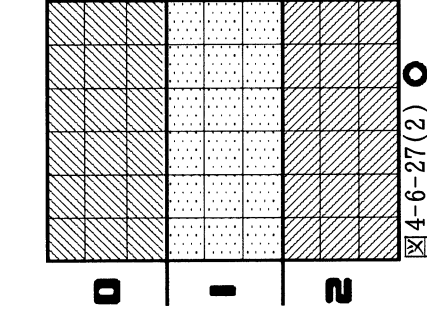
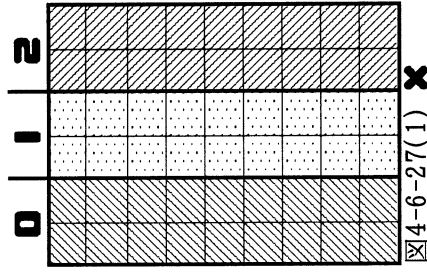


図4-6-26

## ■ パフォーマンス上の考慮点

パイプライン法のパフォーマンス上の考慮点について説明します。

- パイプライン法では1次元方向と2次元方向のどちらでブロック分割することもできますが、境界のデータの通信が必要となるので、計算領域が長方形の場合には以下のように通信量が少なくなる(境界が短くなる)方向(図4-6-27(2), 図4-6-27(4))で分割を行って下さい。

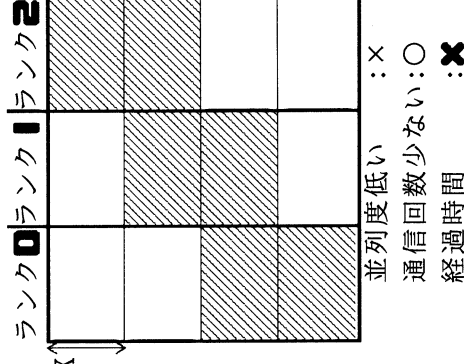
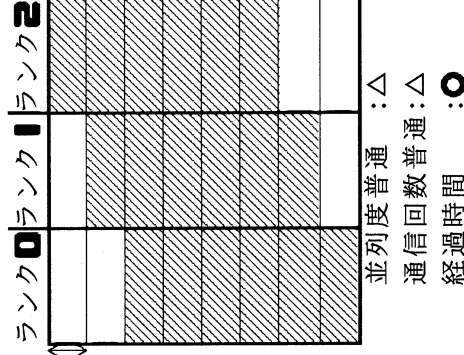
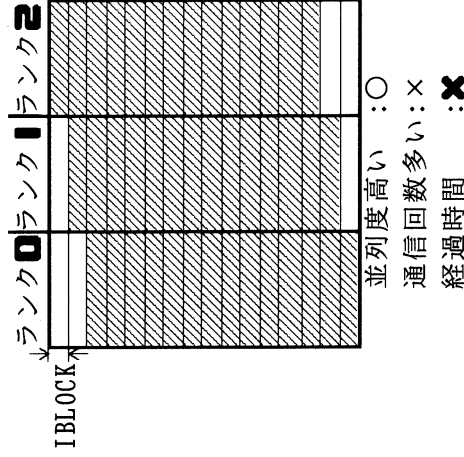


- ブロックの大きさ(IBLOCK)を変えたと、ある値でパフォーマンスがもつともよくなり、それより大きくしても小さくしてもパフォーマンスは低下します。この理由を以下に説明します。

図4-6-28(1)(2)(3)ではブロックの数がそれぞれ16, 8, 4個になっており、この3ケースをまず並列度について比較すると、全プロセスが稼働するのは着色した部分になるので、図4-6-28(1)は並列度が高くなりませんが図4-6-28(3)では低くなります。

次に通信について比較します。各プロセスが行う通信量は境界の長さになるので図4-6-28(1)(2)(3)は同じです。次に通信回数ですが、1つのブロックを処理することに通信が必要になるので、ブロック数の多い図4-6-28(1)では通信回数が多くなり、図4-6-28(3)では少なくなります。4-1節で説明したように、通信を1回行うごとに立上り時間がかかるので、通信回数の点ではブロック数の少ない図4-6-28(3)が最もよくなります。

以上のように、ブロック数は多すぎても少なすぎてもパフォーマンスが低下しますが、最適なIBLOCKの大きさは計算量とも関係するので試行錯誤で決定して下さい。





# 4-6-8 ツイスト分割法 (ADI法などで使用)

本節で説明する方法は、実際のプログラムではそれほど現れることはないのですが、お急ぎの方は飛ばしてもかまいません。

本節では、参考文献[31]に紹介されているツイスト分割法について説明します。ツイスト分割法は、差分法などで使用されるADI法を並列化する場合などに適用することができます。

図4-6-30の①～④のD0ループでは、それぞれ図4-6-31(1)～(4)の計算を行います(図中の「境」は境界の要素を示します)。計算には、矢印に示す依存関係があります。このように1次元方向(↓↑)と2次元方向(→←)の両方向に依存関係のあるパターンを並列化する場合、例えば図4-6-31の点線に示すように、配列Xを2次元方向にブロック分割して①,②のD0ループを普通に並列化し、③,④のD0ループは4-6-7節で説明したパイプライン法で並列化するという方法もありますが、パイプライン法では並列度や通信回数などの問題があります。

このパターンはツイスト分割法で並列化することもできます。ツイスト分割法では、配列Xの1次元目と2次元目をともに「プロセス数」でブロック分割します。例えば図4-6-32(1)に示す10×10の要素を5プロセスでツイスト分割すると、図4-6-32(2)のようになります。図中の着色部分を本書ではブロックと呼びます。図から分かるように、ツイスト分割ではどの行、どの列のブロックも①～④の全プロセスが担当しており、これによって、どの方向の計算も各プロセスのロードバランスのロードバランサーが均等になります。

```

PROGRAM MAIN
PARAMETER(IMAX=6,JMAX=6)
REAL X(0:IMAX+1,0:JMAX+1)
:
D0 J=1,JMAX ① (↓)
  D0 I=1,IMAX
    X(I,J) = X(I,J)+X(I-1,J)
  ENDDO
ENDDO
D0 J=1,JMAX ② (↑)
  D0 I=IMAX,1,-1
    X(I,J) = X(I,J)+X(I+1,J)
  ENDDO
ENDDO
D0 J=1,JMAX ③ (→)
  D0 I=1,IMAX
    X(I,J) = X(I,J-1)+X(I,J)
  ENDDO
ENDDO
D0 J=JMAX,1,-1 ④ (←)
  D0 I=1,IMAX
    X(I,J) = X(I,J+1)+X(I,J)
  ENDDO
ENDDO
:
    
```

図4-6-30

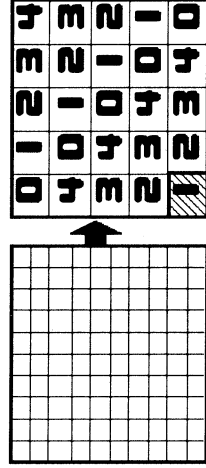


図4-6-32(1)

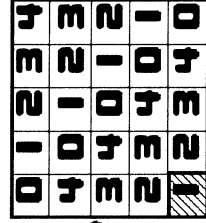


図4-6-32(2)

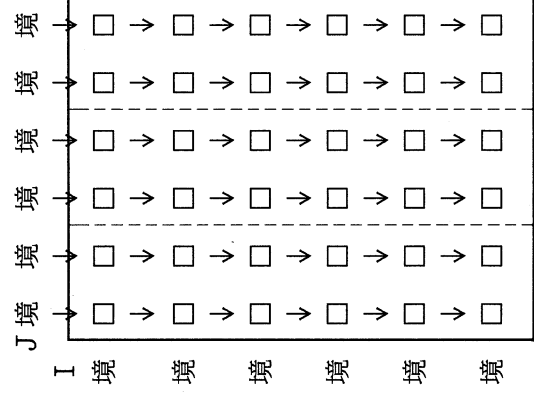


図4-6-31(1) ①のループ

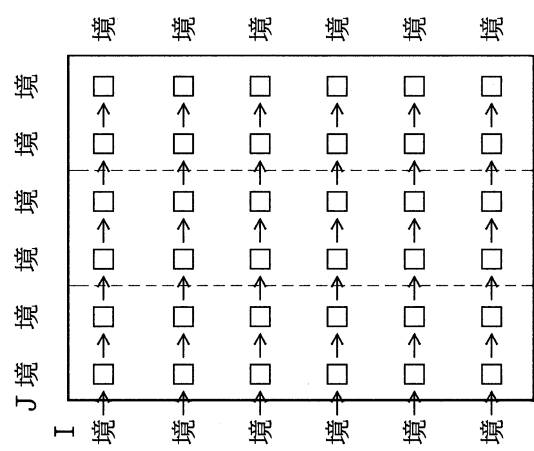


図4-6-31(3) ③のループ

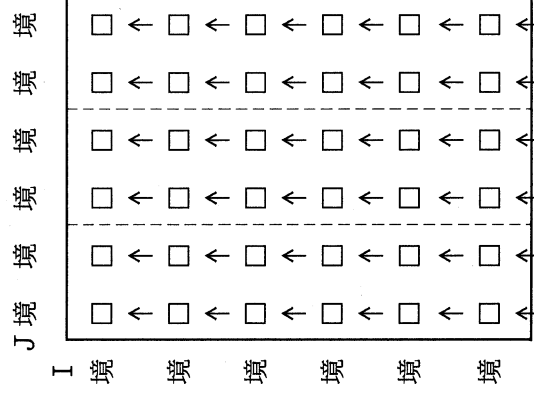


図4-6-31(2) ②のループ

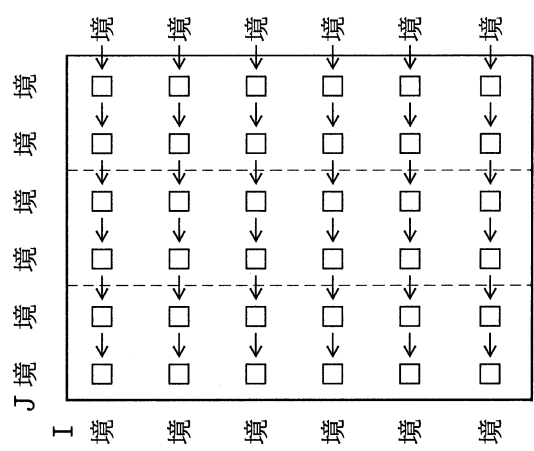


図4-6-31(4) ④のループ

図4-6-33(1)の計算領域を3プロセスでツイスト分割して並列化したプログラムを図4-6-34に、動作を図4-6-35(1)～(4)に示します。図4-6-33(1)の○で囲んだ部分は各ブロックの座標を示します。

- 図4-6-34(1)で、自分のランク値より1つ上と下のランク値をIUPとIDOWNにセットします。この値は0→1→2→0のように循環します。
- (2)で、図4-6-33(1)(2)に示すように、II座標が0,1,2のときに自分のプロセスが担当するブロックのJJ座標を配列JJJにセットします。同様に、JJ座標が0,1,2のときに自分のプロセスが担当するブロックのII座標を配列IIIにセットします。配列IIIとJJJの値は各プロセスで異なることに注意して下さい。
- (3)で、図4-6-33(1)に示すように、各ブロックのI方向とJ方向をブロック分割し、各ブロックの最初と最後の位置を配列ISTA, IEND, JSTA, JENDにセットします。ブロック分割には4-5-4節で紹介したサブルーチンPARA\_RANGEを使用しました。
- (4)と(5)で、データを通信するときに使う派生データ型を、3-5-3-2節で説明した自作のサブルーチンPARA\_TYPE\_BLOCK2を使用して、ブロックごとに作成します(MPI-2が使用できるマシン環境の場合は、MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照)を使用して下さい)。なお、派生データ型は全て配列Xの先頭を起点とします。

例えばブロック(1,1)の場合を図4-6-33(1)の○と□に示します。(4)で、各ブロックのI方向の最初の1行、最後の1行、およびそれぞれの境界の外の1行の派生データ型を作成し、配列ITSTA, ITEND, ITSTAM, ITENDPに設定します。同様に(5)で、J方向の最初の1列、最後の1列、およびそれぞれの境界の外の1列の派生データ型を作成し、JTSTA, JTEND, JTSTAM, JTENDPに設定します。

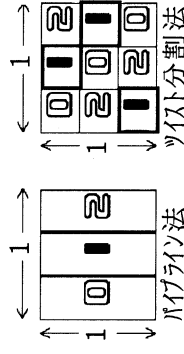
なお、例えば図4-6-33(1)のITSTAM(1,1)と、1つ上のブロックのITEND(0,1)は実際には同一ですが、配列Xを縮小(4-5-7節参照)する場合は配列名が別になっていたので(詳細は省略します)、あえて別の配列にしました。

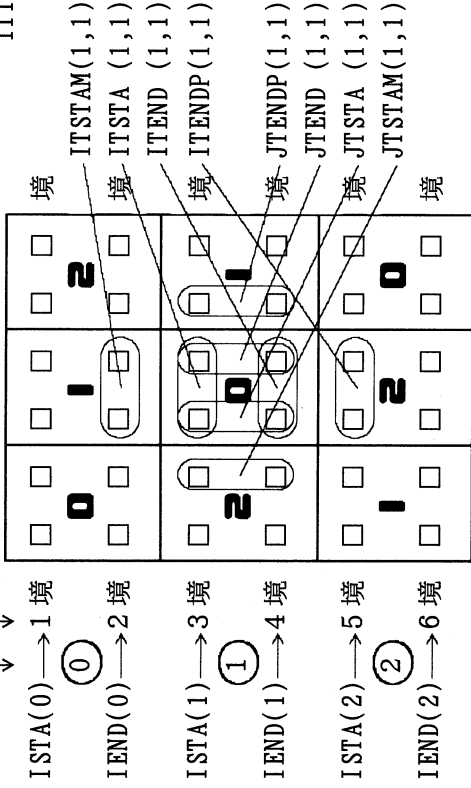
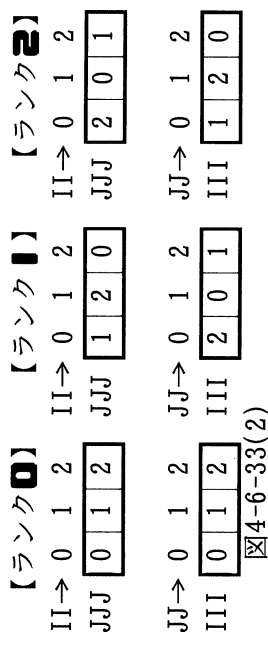
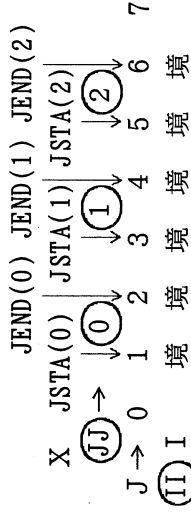
- (6)、(2)、(3)、(4)のD0ループで、図4-6-35(1)～(4)の計算を並列に行います。どの計算も同様なので、ここでは(6)のD0ループについて図4-6-35(1)で説明します。
- (6)のD0ループで、ブロックのI方向の座標IIが0,1,2と変化すると、図4-6-35(1)の一番上、上から2番目、3番目のブロックのうち、自分が担当するブロックのJ方向の座標をJJにセットします。例えばランク0の場合は、II=0のときに自分が担当するブロックのJ方向の座標をJJにセットします。例えばランク0の場合はJJ=0となり、図4-6-35(1)の着色したブロックが選択されます。
- (8)で、各プロセスは自分が担当するブロック内の各要素を計算します。
- 計算終了後、(0)で、計算したブロックの一番下の1行(例えばランク0の場合は図4-6-35(1)のA)を、次に担当するプロセス(本例ではランク2)に送信します。この送信には、各ブロックの一番下の1行を表す派生データ型ITENDを使用します。
- 逆に他のプロセス(本例ではランク1)から送られたデータを(0)で受信します(例えばランク0の場合は図4-6-35(1)のB)。この受信には、各ブロックの一番上の境界の外の1行を表す派生データ型ITSTAMを使用します。

- (6)のループがII=2,3と進み、上から2つ目、3つ目のブロックを同様に計算します。なお、3つ目のブロックの終了後には(0)と(0)の通信は行わないので、(9)のIF文が入っています。

図4-6-31のバターンの場合、パイプライン法(4-6-7節)とツイスト分割法は下記のように長所・短所があり、どちらがよいかは一概には言えません。

- パイプライン法はロードバランス(並列度)の点で問題がありませんが、ツイスト分割法では1,2次元目ともに常に同じ数のプロセスが同じ量の計算を行うので、各プロセスのロードバランスは均等になります。
- パイプライン法では(並列度を高めるために)ブロック数を多くすると通信回数が多くなってしまいました(ツイスト分割法ではブロック数(通信回数)は1,2次元目ともにプロセス数と同じで、(プロセス数が少ない場合は)それほど多くはありません)。
- 右図で例えば3プロセスの場合のランク■の通信量(太線の長さ)を比較すると、  
パイプライン法では $1+1=2$ に対し、ツイスト分割法では $(1/3) \times 4 \times 3=4$ となり、  
ツイスト分割法の方が通信量は2倍多くなります。
- 1つの配列をツイスト分割すると、その配列を使用する全てのループで使われる全ての配列を、動作の複雑なツイスト分割で並列化しなければならず、パイプライン法に比べて修正が面倒になります。





7 境 境 境 境 境 境 境

図4-6-33(1)

```

PROGRAM MAIN
INCLUDE 'mpif.h'
PARAMETER( IMAX=53, JMAX=47, NCPU=10)
INTEGER ISTATUS(MPI_STATUS_SIZE)
REAL X(0:IMAX+1,0:JMAX+1)
INTEGER III(0:NCPU-1), JJJ(0:NCPU-1)
INTEGER ISTA(0:NCPU-1), IEND(0:NCPU-1)
INTEGER JSTA(0:NCPU-1), JEND(0:NCPU-1)
INTEGER ITSTA (0:NCPU-1, 0:NCPU-1)
& , ITEND (0:NCPU-1, 0:NCPU-1)
INTEGER JTSTA (0:NCPU-1, 0:NCPU-1)
& , JTEND (0:NCPU-1, 0:NCPU-1)
INTEGER ITSTAM(0:NCPU-1, 0:NCPU-1)
& , ITENDP(0:NCPU-1, 0:NCPU-1)
INTEGER JTSTAM(0:NCPU-1, 0:NCPU-1)
& , JTENDP(0:NCPU-1, 0:NCPU-1)
REAL XX(0:IMAX+1,0:JMAX+1)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
IUP = MOD(MYRANK+1 , NPROCS) (1)
IDOWN = MOD(MYRANK-1+NPROCS, NPROCS) (1)
DO II=0, NPROCS-1
JJJ(II) = MOD(II+MYRANK , NPROCS)
ENDDO (2)
DO JJ=0, NPROCS-1
III(JJ) = MOD(JJ-MYRANK+NPROCS, NPROCS)
ENDDO

```

```

DO II=0, NPROCS-1
CALL PARA_RANGE
& (1, IMAX, NPROCS, II, ISTA(II), IEND(II))
CALL PARA_RANGE (3)
& (1, JMAX, NPROCS, II, JSTA(II), JEND(II))
ENDDO
DO JJ=0, NPROCS-1
JS = JSTA(JJ)
JE = JEND(JJ)
DO II=0, NPROCS-1
IS = ISTA(II)
IE = IEND(II)
CALL PARA_TYPE_BLOCK2
& (0, IMAX+1, 0, IS, IS, JS, JE,
& MPI_REAL, ITSTA(II, JJ))
CALL PARA_TYPE_BLOCK2
& (0, IMAX+1, 0, IE, IE, JS, JE,
& MPI_REAL, ITEND(II, JJ))
CALL PARA_TYPE_BLOCK2 (4)
& (0, IMAX+1, 0, IS-1, IS-1, JS, JE,
& MPI_REAL, ITSTAM(II, JJ))
CALL PARA_TYPE_BLOCK2
& (0, IMAX+1, 0, IE+1, IE+1, JS, JE,
& MPI_REAL, ITENDP(II, JJ))

```

(次ページへ)

```

CALL PARA_TYPE_BLOCK2
      (0, IMAX+1, 0, IS, IE, JS, JS,
      MPI_REAL, JTSTA(II, JJ))
CALL PARA_TYPE_BLOCK2
      (0, IMAX+1, 0, IS, IE, JE, JE,
      MPI_REAL, JTEND(II, JJ)) (5)
CALL PARA_TYPE_BLOCK2
      (0, IMAX+1, 0, IS, IE, JS-1, JS-1,
      MPI_REAL, JTSTAM(II, JJ))
CALL PARA_TYPE_BLOCK2
      (0, IMAX+1, 0, IS, IE, JE+1, JE+1,
      MPI_REAL, JTENDP(II, JJ))
ENDDO
ENDDO

```

```

DO II=0, NPROCS-1 ( ↓ ) (6)
  JJ = JJJ(II) (7)
  DO J=JSTA(JJ), JEND(JJ)
    DO I=ISTA(II), IEND(II)
      X(I, J) = X(I, J) + X(I-1, J) (8)
    ENDDO
  ENDDO (9)
  IF (II/=NPROCS-1) THEN
    CALL MPI_ISEND
      (X, 1, ITEND(II, JJ) , IDOWN, 1,
      MPI_COMM_WORLD, IREQS, IERR) (10)
    CALL MPI_IRECV
      (X, 1, ITSTAM(II+1, JJJ(II+1)), IUP , 1,
      MPI_COMM_WORLD, IREQR, IERR) (11)
    CALL MPI_WAIT(IREQS, ISTATUS, IERR)
    CALL MPI_WAIT(IREQR, ISTATUS, IERR)
  ENDF
ENDDO

```

```

DO II=NPROCS-1, 0, -1 ( ↑ ) (12)
  JJ = JJJ(II)
  DO J=JSTA(JJ), JEND(JJ)
    DO I=IEND(II), ISTA(II), -1
      X(I, J) = X(I, J) + X(I+1, J)
    ENDDO
  ENDDO
  IF (II/=0) THEN
    CALL MPI_ISEND
      (X, 1, ITSTA(II, JJ) , IUP , 1,
      MPI_COMM_WORLD, IREQS, IERR)
    CALL MPI_IRECV
      (X, 1, ITENDP(II-1, JJJ(II-1)), IDOWN, 1,
      MPI_COMM_WORLD, IREQR, IERR)
    CALL MPI_WAIT(IREQS, ISTATUS, IERR)
    CALL MPI_WAIT(IREQR, ISTATUS, IERR)
  ENDF
ENDDO

```

4-6-34

```

DO JJ=0, NPROCS-1 ( → ) (13)
  II = III(JJ)
  DO J=JSTA(JJ), JEND(JJ)
    DO I=ISTA(II), IEND(II)
      X(I, J) = X(I, J-1) + X(I, J)
    ENDDO
  ENDDO
  IF (JJ/=NPROCS-1) THEN
    CALL MPI_ISEND
      (X, 1, JTEND(II, JJ) , IUP , 1,
      MPI_COMM_WORLD, IREQS, IERR)
    CALL MPI_IRECV
      (X, 1, JTSTAM(III(JJ+1), JJ+1), IDOWN, 1,
      MPI_COMM_WORLD, IREQR, IERR)
    CALL MPI_WAIT(IREQS, ISTATUS, IERR)
    CALL MPI_WAIT(IREQR, ISTATUS, IERR)
  ENDF
ENDDO

```

```

DO JJ=NPROCS-1, 0, -1 ( ← ) (14)
  II = III(JJ)
  DO J=JEND(JJ), JSTA(JJ), -1
    DO I=ISTA(II), IEND(II)
      X(I, J) = X(I, J+1) + X(I, J)
    ENDDO
  ENDDO
  IF (JJ/=0) THEN
    CALL MPI_ISEND
      (X, 1, JTSTA(II, JJ) , IDOWN, 1,
      MPI_COMM_WORLD, IREQS, IERR)
    CALL MPI_IRECV
      (X, 1, JTENDP(III(JJ-1), JJ-1), IUP , 1,
      MPI_COMM_WORLD, IREQR, IERR)
    CALL MPI_WAIT(IREQS, ISTATUS, IERR)
    CALL MPI_WAIT(IREQR, ISTATUS, IERR)
  ENDF
ENDDO

```

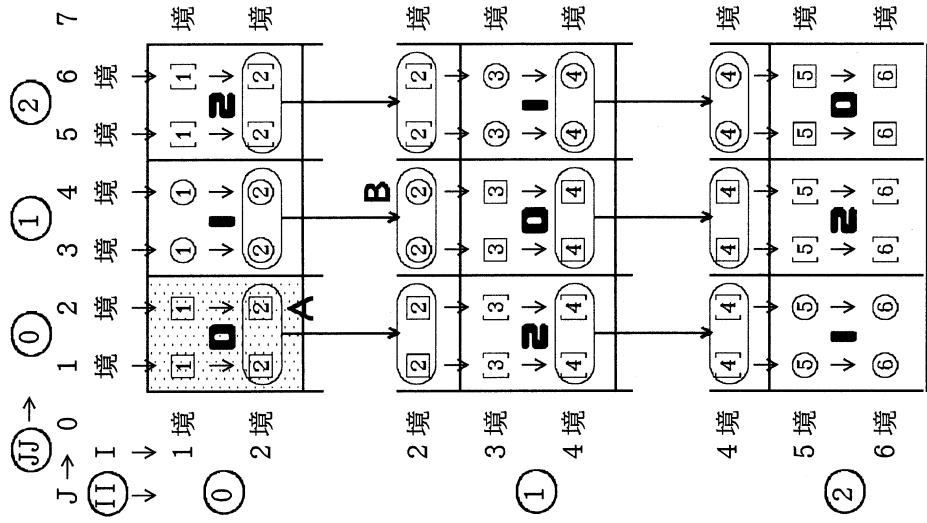


図4-6-35(1) (6)のループ

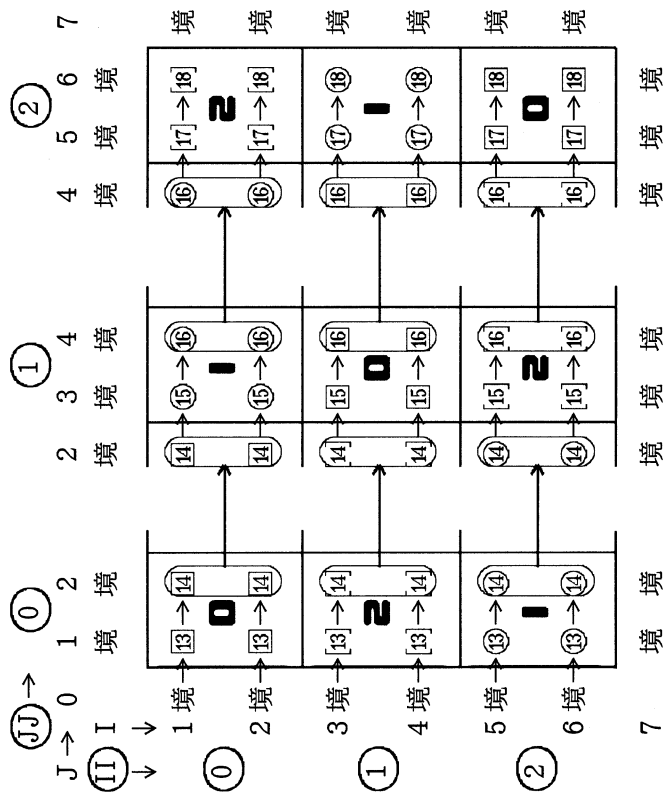


図4-6-35(3) (3)のループ

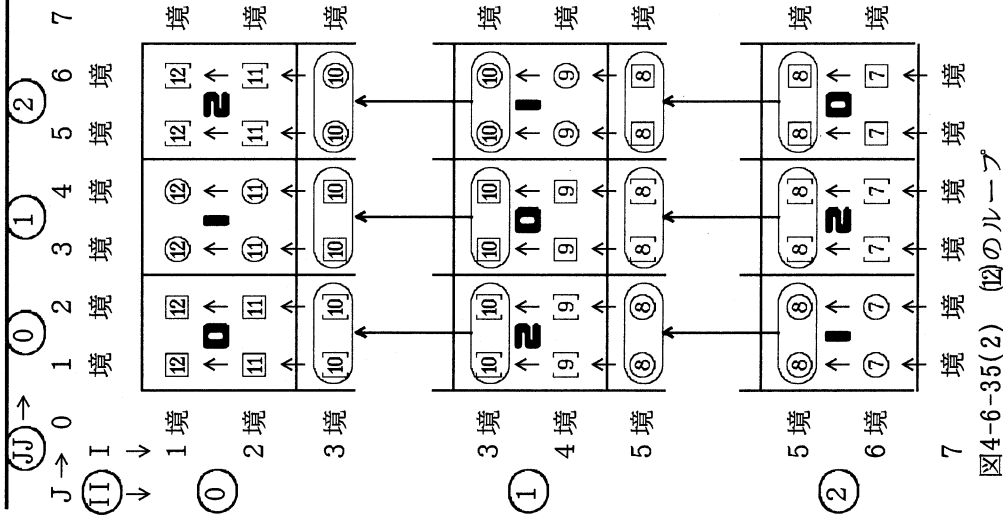


図4-6-35(2) (12)のループ

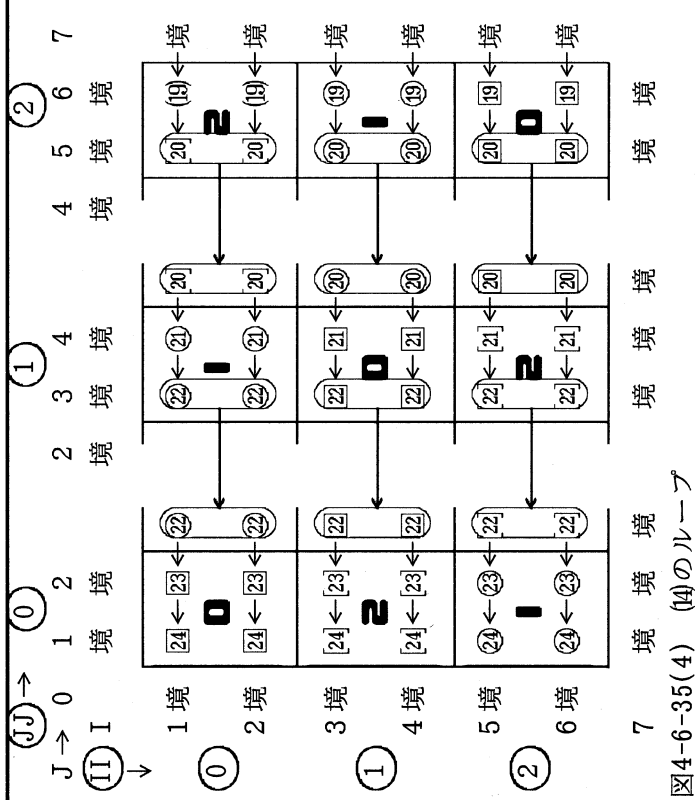


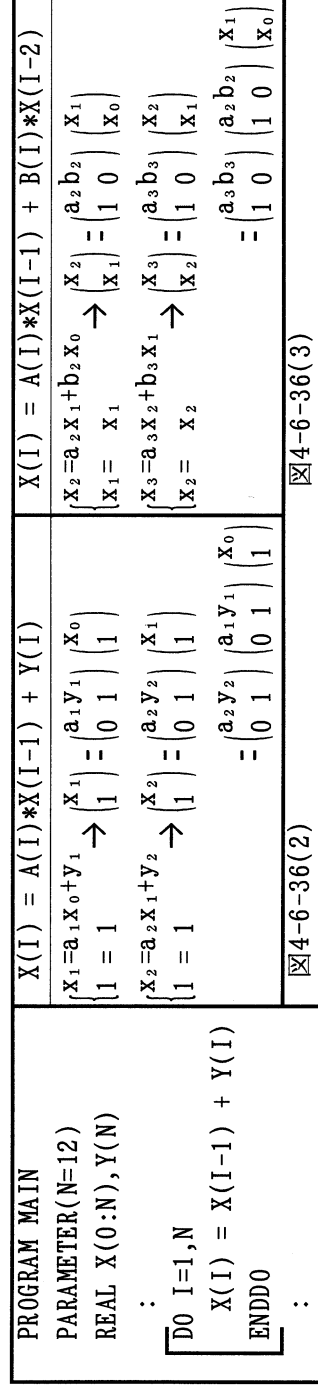
図4-6-35(4) (4)のループ

## 4-6-9 一次逐次演算 (漸化式などで使用)

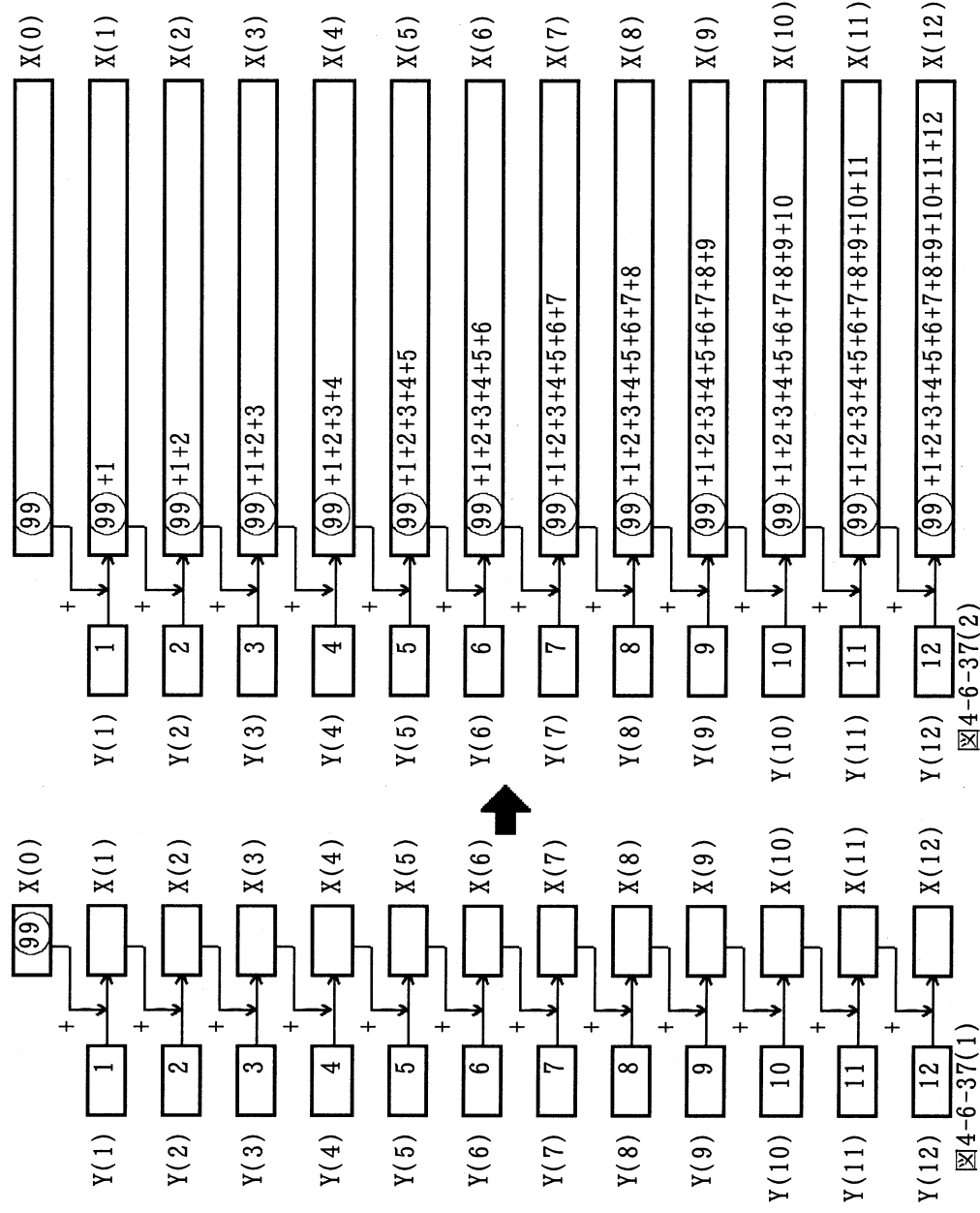
本節で説明する方法は、実際のプログラムではそれほど現れることはないのですが、お急ぎの方は飛ばしてもかまいません。

図4-6-36(1)では、ループの左辺で定義された $X(I)$ が、次のループ反復の右辺で $X(I-1)$ として参照されています。このように、あるループ反復で定義された値がそれ以後のループ反復で参照されることを、「ループの各反復に**帰参照**(または**依存関係**)がある」と言い、通常は**並列化**できません(4-2節参照)。

図4-6-36(1)のように「 $X(I) = X(I-1) + \dots$ 」の形をした、漸化式のような演算を**一次逐次演算**と呼びます(参考文献[18]参照)。図4-6-36(1)には**帰参照**が存在しますが、本節に示すように**並列化**が可能です。図4-6-36(2)(3)も**並列化**が可能です(下記にポイントのみを簡単に記述します(詳細は参考文献[4]参照))。



まず、図4-6-36(1)のループを実行する直前の配列の状態を図4-6-37(1)に、ループを実行した直後の状態を図4-6-37(2)に示します。図4-6-37(1)の配列 $X(0)$ と配列 $Y(1) \sim Y(12)$ には適当な値が入っているとします。図4-6-37(1)(2)から分かるように、 $X(0)$ から $X(1)$ を求め、 $X(1)$ から $X(2)$ を求めるとこの因果関係が連鎖と続くため、通常の方法では**並列化**できないのがお分かりになると思います。



まず並列化したときの動作の概要を説明します。計算する配列が $X(1) \sim X(12)$ 、プロセス数が4とすると、配列 $X$ は図4-6-38(1)のようにブロック分割されます。計算は大きく2つのステップに分かれ、ステップ1で図4-6-38(1)の★(各プロセスが担当する最初の要素)を求め、ステップ2で★から残りの↓を順に求めます。★を求める計算と↓を求める計算は同じなので、各プロセスは同じ計算を2回行うこととなります。

図4-6-36(1)を並列化したプログラムを図4-6-38(2)に、データの動きを図4-6-38(3)に示します。

- 図4-6-38(2)(3)の①で $N(=12)$ 個の要素をブロック分割し、各プロセスが担当する配列 $X$ の下限 $ISTA$ と上限 $IEND$ を、4-5-4節で紹介したサブルーチン`PARA_RANGE`を使用して求めます。
  - ②で変数`SUM`をゼロクリアします。ただしランク0のプロセスだけは $X(0)$ を代入します。
  - ③で、配列 $Y$ のうち自分の担当範囲の値を変数`SUM`に合計します。
  - ④で、MPI-2の集団通信サブルーチン`MPI_EXSCAN`(付録参照)を使用して変数`SUM`の縮約演算を行い、結果を変数`SSUM`に入れます。MPI\_EXSCANを実行すると、図4-6-38(3)の0で囲んだ部分の通信が行われ、自分よりランクが小さい全てのプロセスの`SUM`の合計が自分の`SSUM`に入ります。ただしランク4の`SUM`は使用されず、ランク0の`SSUM`に値は入りません。
- なお、MPI-1の集団通信サブルーチン`MPI_SCAN`を使用しても、ロジックは若干変わりますが、同様の並列化を行うことは可能です。
- ⑤で、`SSUM`に受信データが入っていないランク0のプロセスのみ、 $X(0)$ を`SSUM`に代入します。
  - ⑥で、ランク0~3の各プロセスが担当する最初の要素 $X(1), X(4), X(7), X(10)$ を求めます(これは図4-6-38(1)の★に相当します)。
  - $X(1), X(4), X(7), X(10)$ の値を元に、⑦で残りの要素の計算を行います(図4-6-38(1)の↓に相当します)。

■ 並列化した場合の計算量

この方法では同じ計算(加算)を2度行うため、無駄に感じられるかもしれませんが、例えば要素数が $m$ とすると、単体で実行した場合の加算回数は $m$ 回で、並列化した場合、各プロセスの加算回数は $(2 \times m \text{回}) / (\text{プロセス数})$ となるため(ただし通信時間は除く)、プロセス数が多いときは一応並列化の効果が出ます。

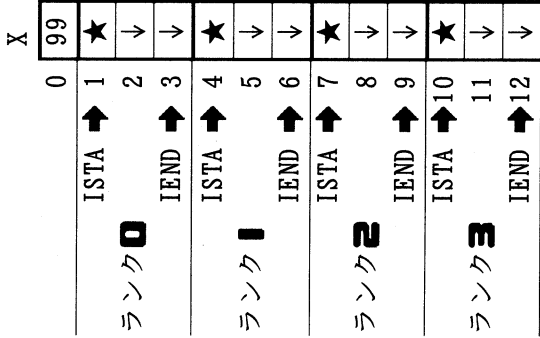


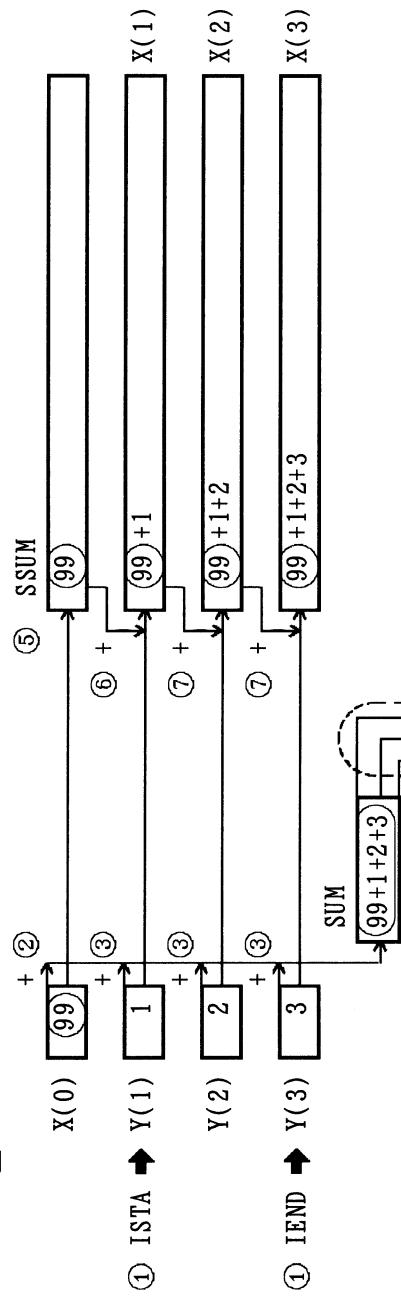
図4-6-38(1)

```

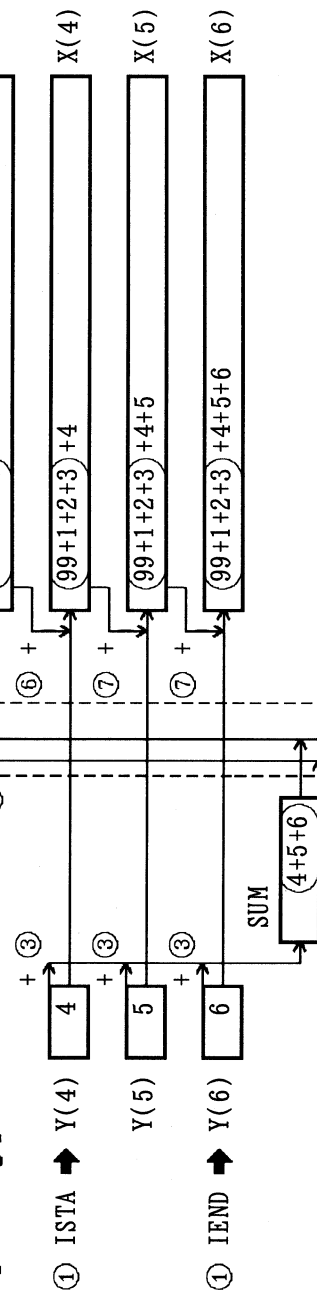
PROGRAM MAIN
  INCLUDE 'mpif.h'
  PARAMETER(N=12)
  REAL X(0:N), Y(N), SUM, SSUM
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  CALL PARA_RANGE(1, N, NPROCS, MYRANK, ISTA, IEND) ①
  :
  SUM = 0.0 ②
  IF (MYRANK==0) SUM = X(0) ②
  DO I=ISTA, IEND ③
    SUM = SUM + Y(I)
  ENDDO
  CALL MPI_EXSCAN(SUM, SSUM, 1, MPI_REAL, ④
    & MPI_SUM, MPI_COMM_WORLD, IERR)
  IF (MYRANK==0) SSUM = X(0) ⑤
  X(ISTA) = Y(ISTA) + SSUM ⑥
  DO I=ISTA+1, IEND ⑦
    X(I) = X(I-1) + Y(I)
  ENDDO
  :
  
```

図4-6-38(2)

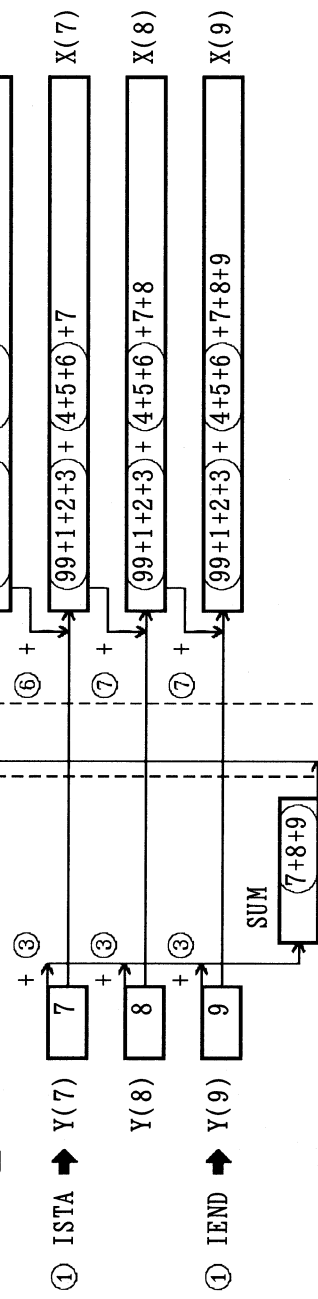
【ランク0】



【ランク1】



【ランク2】



【ランク3】

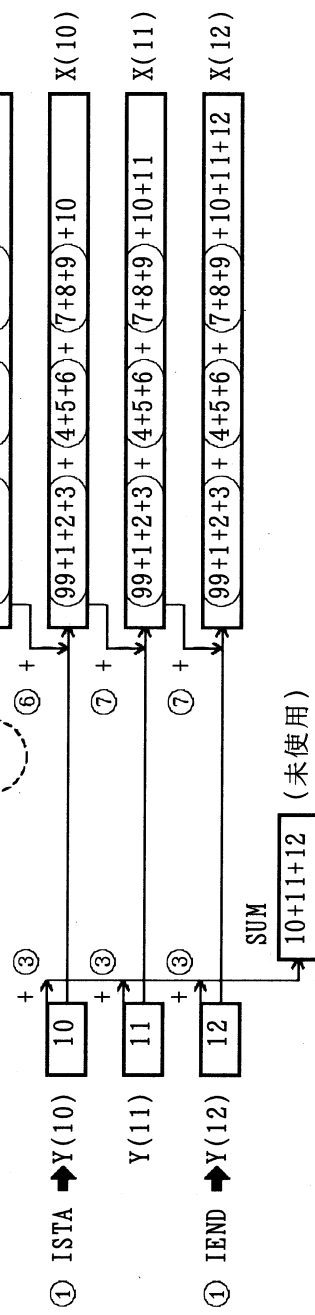


図4-6-38(3)



■ MPI\_ALLTOALLVの場合

図4-6-39(1)では、ランク0~2の各プロセスが、お互いに他のプロセスに個数の異なるデータを送信します。このとき受信側のプロセスは受信する(相手から送られる)個数を知らないとし、この通信は集団通信ルーチンMPI\_ALLTOALLVで行なうことができます。集団通信ルーチンでは、3-3-4節に示したように、送信バッファで指定するバイト数(個数×データ型のバイト数)と受信バッファで指定するバイト数は一致している必要があるため、受信する個数を知る必要があります。本例のように受信する個数を知らない場合の対処方法を図4-6-40のプログラムで説明します。

- 図4-6-40の[3]で、図4-6-39(2)の上段に示すように、ランクiのプロセスに送るデータの個数を配列ISCNT(i)に入れます。何も送らないプロセスには0を入れる必要があるため、[3]の前に[2]でゼロクリアします。
- [4]で、図4-6-39(3)の上段に示すように、ランクiのプロセスに送るデータを送信バッファIRBUF(\*,i)に入れます。
- [5]でMPI\_ALLTOALLの通信を行うと、図4-6-39(2)の下段に示すように、ランクiのプロセスから送られるデータの個数が受信バッファIRCNT(i)に入ります。
- [6]で、図4-6-39(3)の上段に示すISBUFの先頭から、ランクiのプロセス宛てのデータが入る最初の位置までの変位をISDISP(i)に要素数で指定します。本例では、[1]に示すように配列ISBUFの1次元目の大きさが10なので、[6]の下線部で10を指定しました。受信バッファIRBUFについても同様に[6]でIRDISPに指定します。
- [7]でMPI\_ALLTOALLVを使用して通信を行います。このとき、[5]で得られたIRCNTを使用して受信するデータの個数を指定します。その結果、図4-6-39(3)に示す通信が行われます。なお、点線の部分(自分宛ての通信)は、ISCNTとIRCNT内のデータ数がゼロなので通信は行われません。

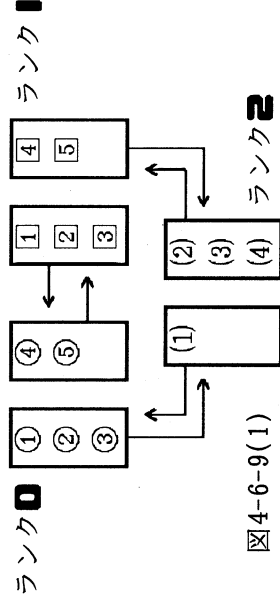


図4-6-9(1)

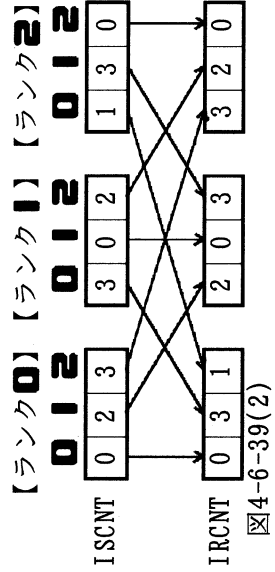


図4-6-39(2)

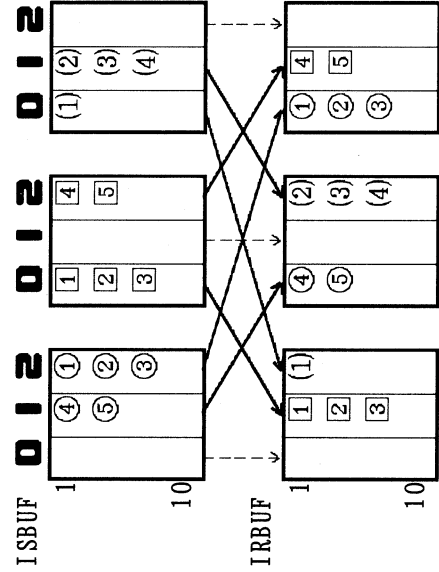


図4-6-39(3)

```

:
PARAMETER(NCPU=3)
INTEGER ISCNT(0:NCPU-1), IRCNT(0:NCPU-1)
INTEGER ISDISP(0:NCPU-1), IRDISP(0:NCPU-1)
REAL SBUF(10,0:NCPU-1), RBUF(10,0:NCPU-1)
:
D0 IRANK=0, NPROCS-1
  ISCNT(IRANK) = 0
ENDDO
ISCNTに送信するデータ数を入れる。
SBUFに送信するデータ数を入れる。
CALL MPI_ALLTOALL(ISCNT,1,MPI_INTEGER,
& IRCNT,1,MPI_INTEGER,
& MPI_COMM_WORLD,IERR)
D0 IRANK=0, NPROCS-1
  ISDISP(IRANK) = IRANK*10
  IRDISP(IRANK) = IRANK*10
ENDDO
CALL MPI_ALLTOALLV(SBUF,ISCNT,ISDISP,MPI_REAL,
& RBUF,IRCNT,IRDISP,MPI_REAL,
& MPI_COMM_WORLD,IERR)
:

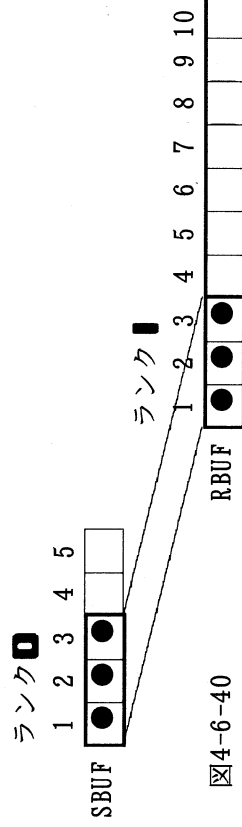
```

図4-6-40

## ■ 1対1通信の場合

1対1通信ルーチンでは、受信バッファで指定するバイト数(データ個数×データ型のバイト数)は、送信バッファで指定するバイト数と等しいか、または大きくても構いません(3-4-5節)。従って、例えばランク1がランク0からデータを受信する場合、受信するデータ個数より受信バッファの方が大きいのが確実であれば、ランク0から送られるデータの個数が分からなくても、1対1通信ルーチンを使用して受信することができます。受信後に、受信したデータの個数を知りたい場合、MPI\_GET\_COUNTを使用します。MPI\_GET\_COUNTの使い方を図4-6-41のプログラムで説明します。

- 図4-6-40に示すように、ランク0の大きさ5の配列SBUFのうち最初の3要素が送信バッファ、ランク1の大きさ10の配列RBUFが受信バッファだします。
- 図4-6-41の[2]で、ランク0は3個のデータをランク1に送信します。
- [3]でランク1は、ランク0からのデータを受信します。このとき受信するデータの個数が分からないので、受信バッファRBUFの大きさである10を指定します。通信の結果、図4-6-40に示すように、受信バッファRBUFの最初の3要素に受信データが入ります。
- ランク1は受信後に、[4]でサブルーチンMPI\_GET\_COUNTを実行します。MPI\_GET\_COUNTの引数には[1]と[3]の二重線に示す配列ISTATUSを指定します(引数の詳細は付録参照)。実行の結果、引数ICOUNTに、[3]で受信したデータの個数(本例では3)が戻り、受信した個数を知ることができます。



```

:
include 'mpif.h'
INTEGER ISTATUS(MPI_STATUS_SIZE) [1]
REAL SBUF(5)
REAL RBUF(10)
:
IF (MYRANK==0) THEN
CALL MPI_SEND(SBUF, 3, MPI_REAL, 1, 1, [2]
& MPI_COMM_WORLD, IERR)
ELSEIF (MYRANK==1) THEN
CALL MPI_RECV(RBUF, 10, MPI_REAL, 0, 1, [3]
& MPI_COMM_WORLD, ISTATUS, IERR)
CALL MPI_GET_COUNT(ISTATUS, MPI_REAL, ICOUNT, IERR) [4]
ENDIF
:

```

図4-6-41

## 4-7 特殊な並列化方法

### 4-7-1 MPMDモデルでの並列化

本節では、複数の異なるプログラムを1つの並列ジョブとして実行するMPMD(Multiple Program Multiple Data)モデルについて説明します(1-2節参照)。

#### 4-7-1-1 実行方法

まずMPMDモデルのプログラムを実行する方法について説明します。

●【方法1】マシン環境によっては、各プロセスごとに異なるロードモジュール名を指定できる機能が提供されている場合があります。これはMPIの機能ではなく、そのマシンの並列処理環境が提供する機能で、提供されているかどうかはマシン環境に依存します。また具体的な指定方法もマシン環境によって異なりますので、マニュアルなどで確認して下さい。

一例を図4-7-2(1)に示します。構造解析と流体解析のロードモジュールを別の名前(a.outとb.out)で作成し、対応表(ファイル)に「ランク0のプロセスはa.out、ランク1のプロセスはb.outを実行する。」と設定します。そして環境変数で対応表のファイル名を指定し、並列ジョブを実行します。

●【方法2】1つのロードモジュール名しか指定できないマシン環境では、通常はMPMDモデルのプログラムを実行することができます。しかし、指定した1つのロードモジュールを各ノードの(異なる)ローカルディレクトリからロードできるマシン環境であれば、MPMDモデルのプログラムを実行することができます。

図4-7-2(2)に示すように、構造解析と流体解析プログラムのロードモジュール名を例えばa.outのように同一名とした(ただし中身は異なります)、各ノードのローカルディレクトリに置いて並列に実行します。すると各ノードで異なる自身のa.outがロードされるので、MPMDモデルのプログラムを実行することができます。

●【方法3】図4-7-2(3)に示すように、構造解析と流体解析のプログラムのメインルーチンをサブルーチン化し、その上に形式的なメインルーチンを置きます(最近流行?の××ホールディングス(持株会社)のようになります)。そしてメインルーチンの最初にランクの値を調べ、いずれかのサブルーチンをコールします。この方法は、2つのプログラムで同じ名前のCOMMONを使っていると、うまくいきません。

●【方法4】MPI-1の拡張機能であるMPI-2(3-1節参照)には、「プロセスの生成と制御」という機能があり、並列ジョブの実行中に別プロセスを作成することができますが、本書では説明を省略します。

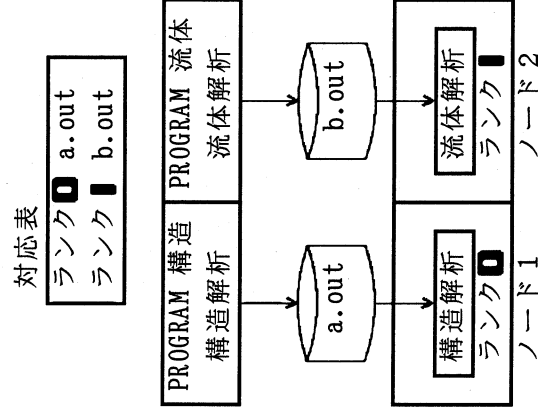


図4-7-2(1)

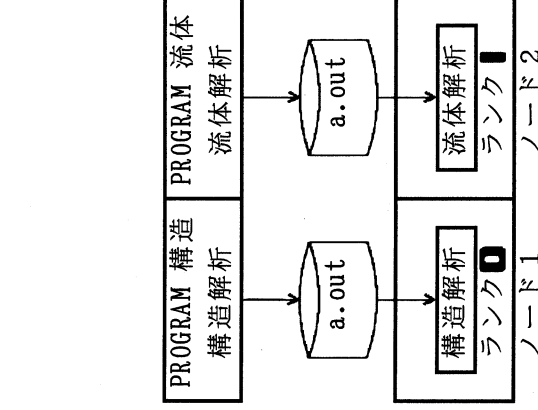


図4-7-2(2)

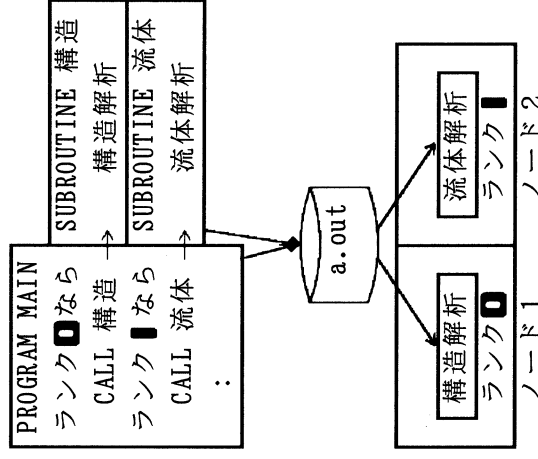


図4-7-2(3)

## 4-7-1-2 連成解析

連成解析のプログラム例を示します。図4-7-3(1)は構造解析、図4-7-3(2)は流体解析のプログラムで、これらを並列に実行します。両プログラムは似ているように見えますが、①に示すように使用する配列は異なっており、全く異なるプログラムだと考えて下さい。

両プログラムとも②に示す同じタイムステップ(変数名はわざと変えています)で反復し、毎回の反復でまず③の計算を行います。

計算が終了すると、④～⑦でお互いに計算結果(圧力や温度などの物理量)を交換します。このとき、例えばランク0は④で配列Aのデータを送信し、ランク1は⑤でデータを(配列Aではなく)配列AAに受信しますが、これは全く問題ありません。④のMPI\_SENDは単に『配列Aのデータをランク1に送れ』と指示するだけで、ランク1が自分と同じプログラムなのか、あるいは受信する配列名は何なのかは全く関知しません。同様に⑤のMPI\_RECVは単に『ランク0から送信されたデータを受信して配列AAに入れよ』と指示するだけで、ランク0が自分と同じプログラムなのか、あるいは送信元の配列名は何なのかは全く関知しません。

⑧でランク0は、ランク1から送られた配列Xの値を使用して配列Aの値を補正します。ランク1も同様です。

補足ですが、④～⑦で、ランク0とランク1ではMPI\_SENDとMPI\_RECVの順序が逆になっています。これは3-4-4節で説明したデッドロックの発生を抑止するためです。もちろんMPI\_ISEND, MPI\_IRecv, MPI\_WAITを使用しても構いません。

### ランク0

```

PROGRAM KOUZOU
PARAMETER (N=1000)
REAL A(N),X(N)
DO I=1,N
  A(I) = ~      ③ 構造解析の計算
ENDDO
CALL MPI_SEND(A,N,MPI_REAL,1,1, ④
& MPI_COMM_WORLD,IERR)
CALL MPI_RECV(X,1,MPI_REAL,1,1, ⑥
& MPI_COMM_WORLD,ISTATUS,IERR)
DO I=1,N
  A(I) = A(I)+X(I)  ⑧ 流体解析から
ENDDO 送られたデータを使って補正
ENDDO

```

図4-7-3(1) 構造解析

### ランク1

```

PROGRAM RYUTAI
PARAMETER (NN=1000)
REAL XX(NN),AA(NN)
DO ISTEP = 1, 100 ②
DO I=1,NN
  XX(I) = ~      ③ 流体解析の計算
ENDDO
CALL MPI_RECV(AA,NN,MPI_REAL,0,1, ⑤
& MPI_COMM_WORLD,ISTATUS,IERR)
CALL MPI_SEND(XX,NN,MPI_REAL,0,1 ⑦
& MPI_COMM_WORLD,IERR)
DO I=1,NN
  XX(I) = XX(I)*AA(I) ⑧ 構造解析から
ENDDO 送られたデータを使って補正
ENDDO

```

図4-7-3(2) 流体解析

上記の例では各プログラム内は並列化されていませんが、並列化することも可能です。図4-7-4では、構造解析のプログラムをランク0～2で、流体解析のプログラムをランク3～5で並列に実行しています。この場合、3-6節で説明したように、構造解析と流体解析でそれぞれ新グループを作成し、新グループ内のプロセス間の通信は新グループのコミュニケータ(以下の例ではIKOUZOUとIRYUTAI)を指定して行い、構造プロセスと流体プロセス間の通信はMPI\_COMM\_WORLDを指定して行います。

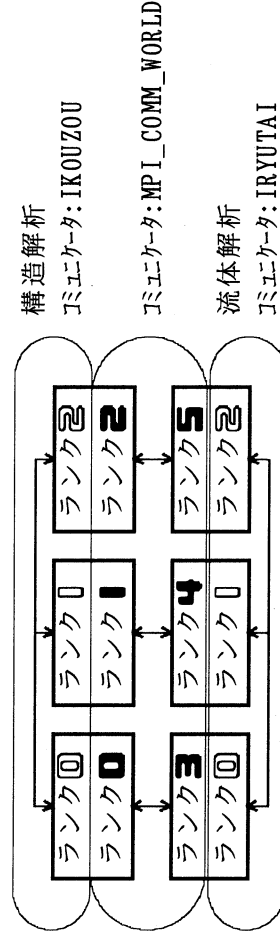


図4-7-4

### 4-7-1-3 マスター・スレーブ方式

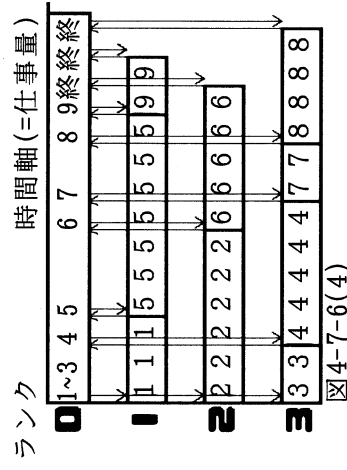
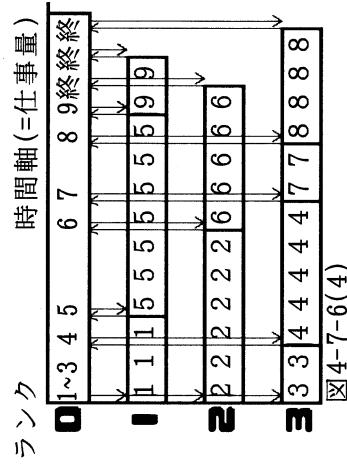
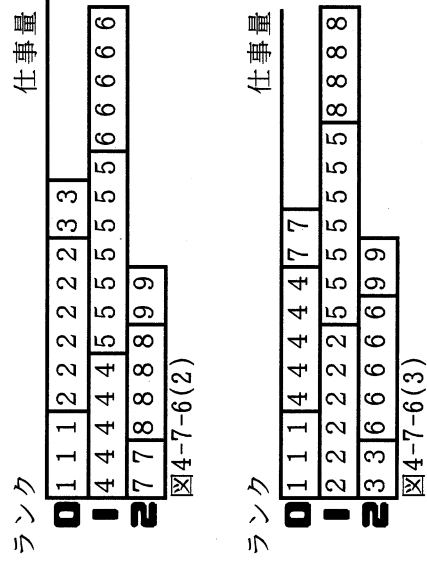
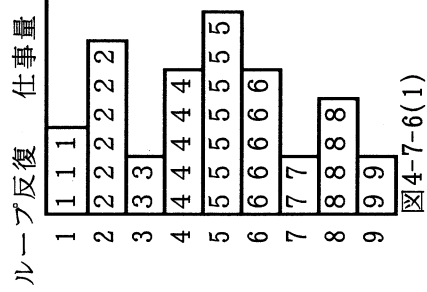
#### ■ マスター・スレーブ方式とは

本節では、1-2節で紹介したマスター・スレーブ(またはマスター・ワーカー)方式について説明します。マスター・スレーブ方式では、1つの親プロセス(係長)が、その全ての子プロセス(係員)に仕事を与え、子プロセス(係員)は仕事が終了したら親プロセス(係長)に報告し、親プロセス(係長)はその子プロセス(係員)に次の仕事を渡す、という動作をします。

マスター・スレーブ方式で並列化するのはどのような場合でしょうか？ 例えば反復回数が9回のD0ループがあり、ループの各反復での仕事量が図4-7-6(1)のように著しくバラついているとします。これを、マシン性能(処理能力)が同一である3つのノードで並列に実行する場合、図4-7-6(2)のようにブロック分割しても、図4-7-6(3)のようにサイクリック分割しても、一般にロードバランスは不均等になります。

これをマスター・スレーブ方式で並列化した場合の動作を図4-7-6(4)に示します。まずランク0(親プロセス)はランク1~3(子プロセス)に1~3反復目の仕事を与え、各子プロセスは与えられた反復での仕事を行います。最も仕事量の少ないランク3が最初に終了し、ランク0に仕事の終了を報告します。ランク0はその時点での次の仕事である4反復目の仕事をランク3に与えます。以下同様に計算が行われます。

このようにマスター・スレーブ方式では、子プロセスが仕事を終了するたびに親プロセスはその時点での次の仕事を与え、各子プロセスは絶え間なく仕事をするので、ロードバランスは(通常)ほぼ均等になります。すなわちマスター・スレーブ方式は、ループの各反復での仕事量が著しくバラついている場合に有効と



上記は、各仕事量が著しくバラついている各マシン性能(処理能力)は同一の場合でしたが、逆に、ループの各反復での仕事量は一定で、マシン性能(処理能力)がバラついている場合も有効となります。例えば性能の異なるノードから構成される並列計算機で並列ジョブを実行する場合や、性能の同じノードであっても、一部のノードで他のプログラムが動いている、ノードによって処理能力が異なっている場合などが考えられます。

図4-7-7(1)に示すように、反復回数が9回のD0ループがあり、各反復での仕事量は同じであるとしても、またランク1を実行するノードはマシン性能が高く、ランク2のノードは普通、ランク3のノードは低いとします。これをマスター・スレーブ方式で並列化した場合の動作を図4-7-7(2)に示します。

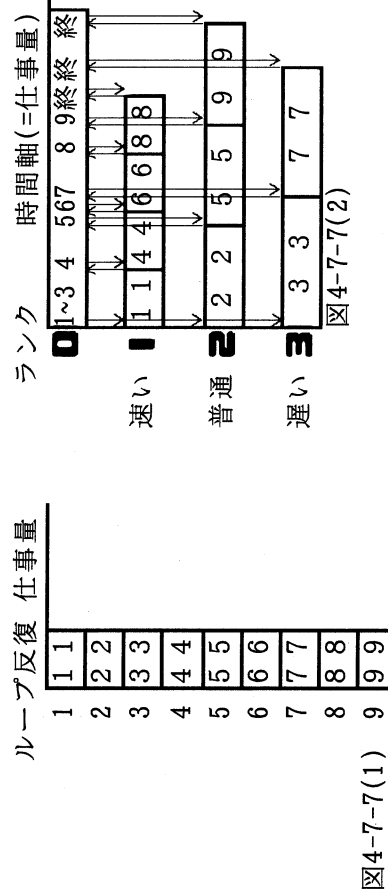
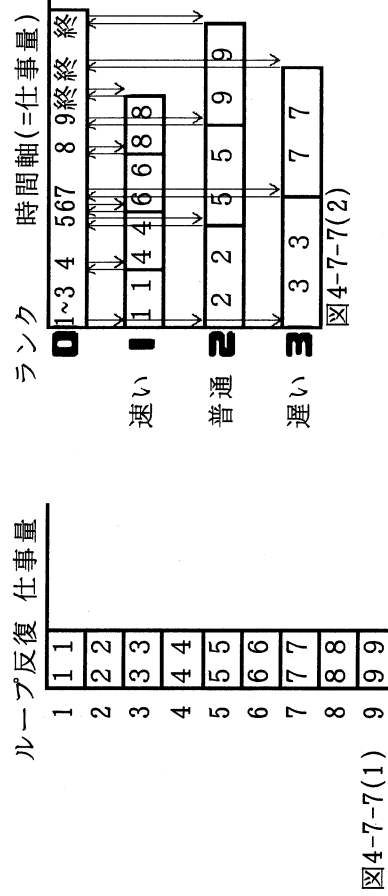


図4-7-7(1)

図4-7-7(2)

マスター・スレーブ方式に関していくつか補足します。

● 前述の例では、通信は親プロセスと子プロセス間のみで、子プロセス間の通信はありませんでしたが、子プロセス間で通信が必要なプログラムの場合、その時点で各子プロセスがどのループ反復を実行しているかはタイミングによって変わるので、マスター・スレーブ方式で並列化するのは難しくなります。

● 前述の例では、親プロセスは仕事の割り振りのみをしているので、通常、子プロセスと比べて仕事量は少なくなりますが、その場合、どれか1つのノードに親プロセスと1つの子プロセスの両方を実行させるという方法もあります。

#### ■ マスター・スレーブ方式のプログラム例

図4-7-8の単体プログラムを例に、マスター・スレーブ方式で並列化する方法を説明します。①のループでは1回の反復で1つの仕事を行います。

各反復の②で[0~1の一樣乱数]を1つ発生させ、それを10倍して整数NSLEEPに代入します。つまりNSLEEPには0~9のいずれかの整数が入ります。次に③でNSLEEP秒間スリープします。従ってループの各反復での経過時間は、0~9秒のいずれかにバラつくこととなります。なお、[0~1の一樣乱数]を生成する組み込み関数、およびスリープさせる組み込み関数はコンパイラによって異なるので、ここでは具体的に書きませんでした。

```
PROGRAM MAIN
DO I=1,10
  NSLEEP = [0~1の一樣乱数]*10.0
  NSLEEP秒間スリープする。
  PRINT *, 'SLEEP TIME = ',NSLEEP
ENDDO
END
```

図4-7-8

図4-7-8をマスター・スレーブ方式で並列化したプログラムのマスター部分を図4-7-9(1)に、スレーブ部分を図4-7-9(2)に示します。なお、本例ではマスター部分とスレーブ部分を同一プログラムとしましたが、別プログラムにする場合は(4-7-1-1節参照)、図4-7-9(1)(2)の②,④,⑥を除去し、①を③にコピーし、④にEND文を追加して下さい。

#### (1) マスタープログラムの動作

- まず図4-7-9(1)の④で、任意の子プロセスから仕事の要求が来るのを待ちます。送信元のプロセスが「任意」なので、下線に示すようにワイルドカードMPI\_ANY\_SOURCEを指定します(注)。④は単に子プロセスからの仕事の要求を受けるだけで、受信するデータそのものには意味がないので、受信バッファーには適当な変数(IDUMMY)を指定します。
- ④が完了すると、配列ISTATUSには送信元のプロセスのランクがセットされるので、⑤でランクの値を取り出して変数IDESTにセットします(注)。
- ⑥で、乱数を発生させて子プロセスをスリープさせる秒数NSLEEPにセットし、⑦でランクがIDEST(つまり④の受信の送信元)のプロセスに対してNSLEEPの値を送信します。
- ③のD0ループを10回反復させた後、⑧のD0ループで全子プロセスを終了させます。まず⑨で任意の子プロセスから仕事の要求が来たら、仕事の終了を知らせるために、そのプロセスに対して⑩でNSLEEPに『-999』をセットし、⑪で送信します。⑧のD0ループは子プロセスの個数分(=全プロセス数-1)だけ反復します。

(注)「付録」のサブルーチンMPI\_RECV, MPI\_I\_RECV, MPI\_WAITの引数を参照して下さい。

(2) スレーブプログラムの動作

●各子プロセスは、図4-7-9(2)の⑮で無限ループします(このD0文はFortran90でサポートされている機能で、無限ループを意味します)。まず⑯で親プロセスに仕事の要求をします。これは⑭または⑰に対応する送信になります。

●⑰で親プロセスからの応答を受信します。⑰から来た応答の場合はスリープする秒数がNSLEEPにセットされ、⑱から来た応答の場合はジョブの終了を示す『-999』がNSLEEPにセットされます。

●⑲でNSLEEPの値を調べ、『-999』の場合は⑲で終了します。『-999』以外の場合は⑳でNSLEEP秒間スリープします。

4プロセス(ランク0)が親プロセス、ランク1~3が子プロセス)で実行した場合の結果を右に示します。

```
[aoyama@node01]~/u/aoyama: mpirun -np 4 a.out
MYRANK = 3 SLEEP TIME = 3
MYRANK = 3 SLEEP TIME = 1
MYRANK = 1 SLEEP TIME = 5
MYRANK = 2 SLEEP TIME = 5
MYRANK = 2 SLEEP TIME = 5
MYRANK = 1 SLEEP TIME = 6
MYRANK = 3 SLEEP TIME = 7
MYRANK = 3 SLEEP TIME = 1
MYRANK = 2 SLEEP TIME = 4
MYRANK = 1 SLEEP TIME = 8
```

```
PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
IF(MYRANK==0) THEN マスタースレーブプロセスの処理②
DO I=1, 10 本処理 ③
CALL MPI_RECV(IDUMMY, 1, MPI_INTEGER,
MPI_ANY_SOURCE, 1, MPI_COMM_WORLD,
ISTATUS, IERR) ④
IDEST = ISTATUS(MPI_SOURCE) ⑤
NSLEEP = [0~1の一樣乱数]*10.0 ⑥
CALL MPI_SEND(NSLEEP, 1, MPI_INTEGER,
IDEST, 1, MPI_COMM_WORLD, IERR) ⑦
ENDDO
DO I=1, NPROCS-1 終了処理 ⑧
CALL MPI_RECV(IDUMMY, 1, MPI_INTEGER,
MPI_ANY_SOURCE, 1, MPI_COMM_WORLD,
ISTATUS, IERR) ⑨
IDEST = ISTATUS(MPI_SOURCE)
NSLEEP = -999 ⑩
CALL MPI_SEND(NSLEEP, 1, MPI_INTEGER,
IDEST, 1, MPI_COMM_WORLD, IERR) ⑪
ENDDO
CALL MPI_FINALIZE(IERR) ⑫
```

図4-7-9(1)

```
ELSE
DO
CALL MPI_SEND(IDUMMY, 1, MPI_INTEGER, 0, 1,
MPI_COMM_WORLD, IERR) ⑮
CALL MPI_RECV(NSLEEP, 1, MPI_INTEGER, 0, 1,
MPI_COMM_WORLD, ISTATUS, IERR) ⑰
IF(NSLEEP== -999) THEN
CALL MPI_FINALIZE(IERR) ⑲
STOP
ENDIF
NSLEEP秒間スリープする。 ⑳
PRINT *, 'MYRANK = ', MYRANK,
'SLEEP TIME = ', NSLEEP
ENDDO
ENDIF
END
```

図4-7-9(2)

## 4-7-2 フラットMPIとハイブリッド並列

本節では、図4-7-10(1)に示すように各ノードに複数CPUが接続されていて、ノード間には分散メモリー型、ノード内には共有メモリー型(以後SMP)となっている並列計算機での並列化について説明します。なお、具体的なコンパイル・実行方法はマシン環境に依存しますので、マニュアルなどを参照して下さい。

### ■ 2CPUを使用する場合

元のプログラムで、図4-7-10(1)のD0ループは1~100まで反復するとします。これを並列化し、4CPUのうち2CPUを利用して実行する3つの方法を比較します。どの方法が最もパフォーマンスがよいかは、プログラムの特性に依存し、ケースバイケースです。

- 図4-7-10(1)では、MPIで並列化したプログラムを、ノード間にまたがる2CPUを使用して並列に実行しています。また、 $\hookrightarrow$ の数字は、各プロセスが担当するD0ループの下限と上限を示します。 $\leftrightarrow$ に示すように、ノード間の通信のオーバーヘッドがかかります。
- 図4-7-10(2)では、SMPで並列化したプログラムを、ノード内の2CPUを使用して並列に実行しています。SMPの並列化では基本的にD0ループが並列化の対象となり、①に示すように指示行を指定して強制的に並列化するか、またはコンパイラが自動的に並列化します。並列化されたD0ループは、スレッド(プロセスのよくなもの)に分かれて実行が行なわれますが、本書ではこれ以上の説明は省略します。SMPでは並列化されたD0ループの前後(D0ループに入る時と終了するとき)でオーバーヘッドが発生します。また同一ノード内の複数CPUを使用するため、資源の競合(メモリーやCPUなど)によるオーバーヘッドが発生します。
- 図4-7-10(3)では、MPIで並列化したプログラムを、同一ノード内の2CPUを使用して並列に実行しています。同一ノード内のCPU間でMPIの通信を行う場合、マシン環境によっては、実際に通信せずにメモリーコピーで代替し、通信時間を短縮できる場合があります。一方、同一ノード内の2CPUを使用するため、図4-7-10(2)と同様に資源の競合(メモリーやCPUなど)によるオーバーヘッドが発生します。

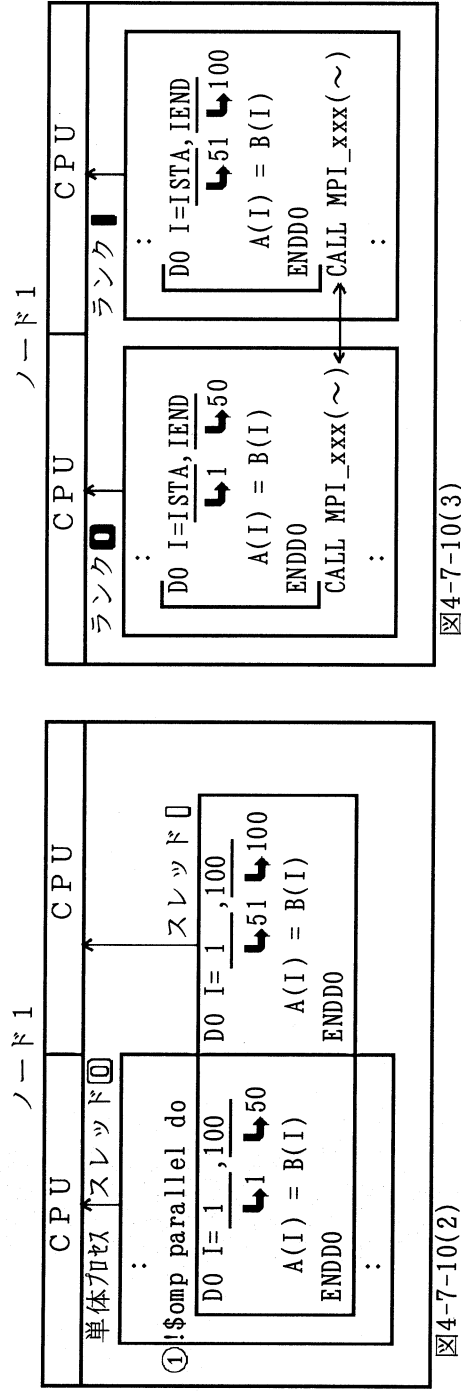
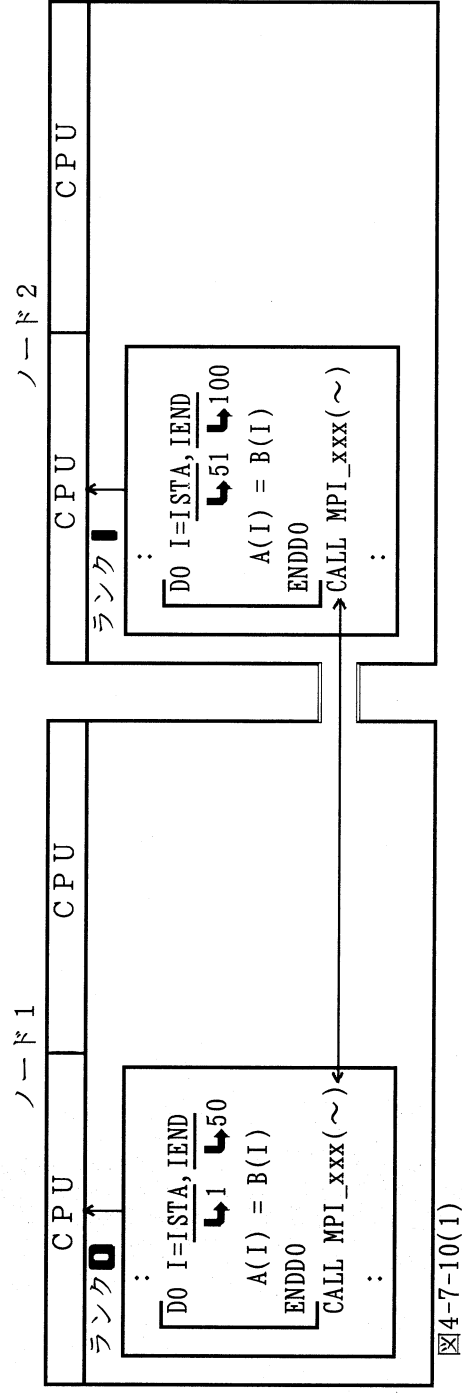


図4-7-10(3)



■ 4CPUを使用する場合

図4-7-10(1)の4CPU全てを利用して並列化する2つの方法と比較します。どちらの方法がパフォーマンスがよいかはプログラムの特性に依存し、ケースバイケースです。

● 図4-7-10(4)では、MPIで並列化したプログラムを4プロセスで実行しています。このようにMPIのみで並列化する方法をフラットMPIと呼びます。ノード間の通信によるオーバーヘッドと、ノード内の資源の競合(メモリーやCPUなど)によるオーバーヘッドが発生します。

● 図4-7-10(5)では、MPIで並列化したプログラムを2プロセスで実行し、さらにDOループをSMPで並列化しています。このようにMPIとSMPの両方で並列化する方法をハイブリッド並列と呼びます。ノード間の通信によるオーバーヘッド、図4-7-10(2)と同様にスレッド起動によるオーバーヘッドが発生します。

● 配列A,Bを縮小(4-5-7節参照)しなかった場合に、並列ジョブ全体で使用できる配列の大きさを比較してみます。配列Aのうちランク0~3のプロセスおよびスレッド0,1が担当する部分を、図4-7-10(4)(5)の下の着色した部分に示します。図4-7-10(4)では配列Aのうち25%しか使用していませんが、図4-7-10(5)では50%を使用しています。以上から次の事が分かります。MPIで並列化した図4-7-10(4)のプログラムに対し、大量のメモリーを使用するため配列の縮小を試み、何らかの理由(ロジックの修正が面倒など)で縮小できなかった場合、図4-7-10(5)のようにSMPでさらに並列化すれば、図4-7-10(4)よりもメモリーを多く使用することができま

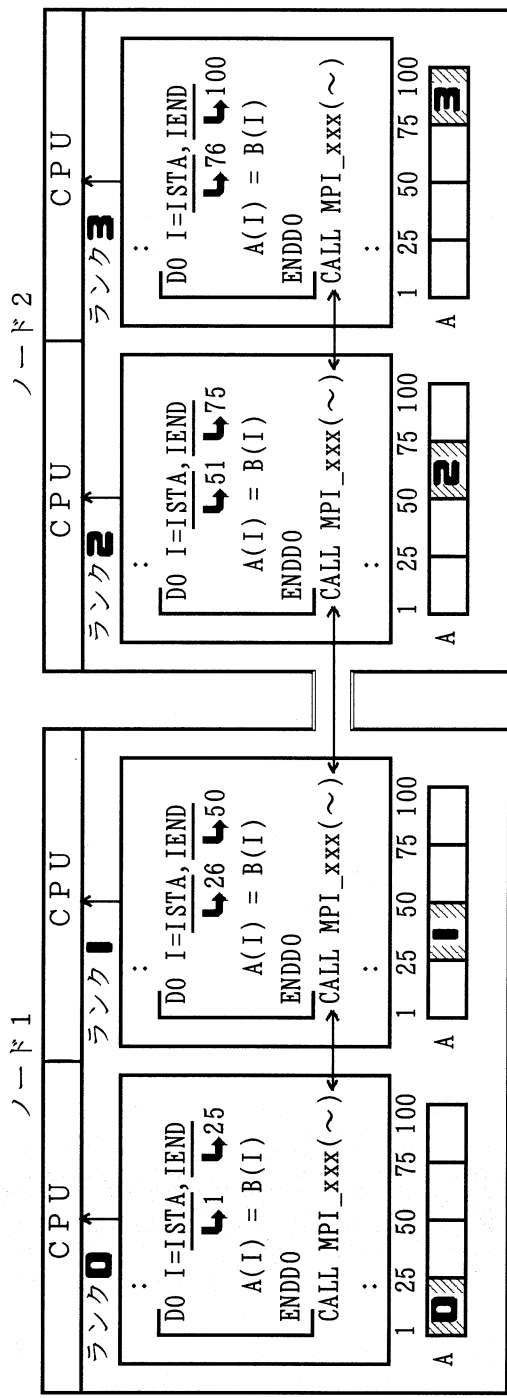


図4-7-10(4)

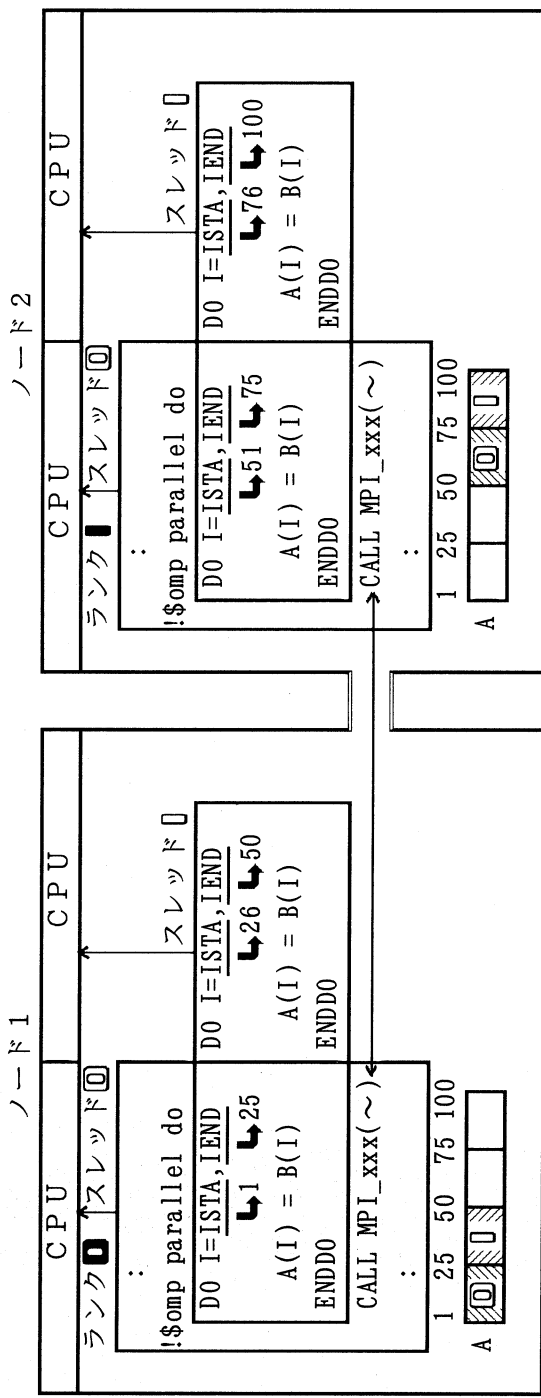


図4-7-10(5)

## 4-8 並列化の手順と考慮点

本章ではプログラムを実際に並列化するための方法について述べてきました。本節ではそのまとめとして、並列化を行う手順およびパフォーマンスの評価方法について説明します。

### 4-8-1 並列化の手順

プログラムの並列化には、単体で稼働しているプログラムを並列化する場合と、最初から並列化を想定してプログラムを作成する場合があります。ここでは前者を想定し、単体プログラムを並列化する一般的な手順について説明します。

#### (1) ホットスポットの特定

まず、単体でプログラムを実行した場合のホットスポット(CPU時間を多く費やしているサブルーチンまたはループ)を、profコマンドなど(マシン環境に依存します)を使用して調べます。

#### (2) 単体チューニング

ホットスポットの部分について単体チューニング(参考文献[6]参照)を行います。単体チューニングを行なって効果が出た場合には、CPU時間の分布が変わるので、チューニングが終了したら再度profコマンドなどでサブルーチンごとのCPU時間の割合を調べます。

#### (3) 並列化の全体方針の決定

profコマンドなどの報告書をもとに、4-3節のどのパターンで並列化するかを決定します。なお、並列化を行う前に、以下の方法で並列化の効果を見積り、並列化作業自体を行うべきかどうかを検討することができます。

● 並列化する部分のCPU時間の割合がわかれば、並列化した場合の効果の上限を4-1節に述べたアムダールの法則を使って知ることができます。例えば並列化する部分が全体の8割しか占めていなければ、理想的な並列計算機を∞台使用しても速度向上率は5倍以上にはなりません。

● 例えば並列化した部分のCPU時間が400秒かかり、並列化によって発生する通信量が大体nバイトであることがわかっているとしめます。簡単なプログラムでnバイトの通信時間を測定し、例えば4ノードで10秒かかったすると、CPU時間は4ノードではば1/4の100秒になりますので、 $400 \div (100+10) = 3.64$ 倍の速度向上率になると見積もることができます。

#### (4) 並列化の詳細方針の決定

以下の点について並列化の詳細方針を決定します。

- 多重ループの場合、どのループを分割するか。言いかえると、多次元配列をどの次元で分割するか。また分割方法をブロック分割、サイクリック分割、ブロック・サイクリック分割のどれにするか。
- 配列の縮小を行うかどうか。
- どの変数(配列)をどこでどのプロセスに通信すれば、並列ジョブとして結果が正しくなり、しかも通信量が最小になるか。
- 並列版の数値計算ライブラリーで提供されているサブルーチンを使用するかどうか。

(5) 並列に動作させる最低限の修正

詳細方針が決定したら実際の並列化に入ります。まず、並列に動作させるための最低限の修正を行います。並列化する前のプログラムを図4-8-1(1)に示します。サブルーチンSUBの中に並列化したいループ(⑦)があり、その反復回数Nが④で与えられているとします。

並列に動作させるための最低限の修正を加えたプログラムを図4-8-1(2)に示します。

● 図4-8-1(2)のプログラムの開始直後、⑩で「三種の神器」を実行して全プロセス数NPROCS、自分のランク値MYRANKを取得します。

これらの変数は、並列化する各サブルーチン(本例ではSUBのみ)で使用しますが、引数で渡すとプログラムがゴチャゴチャするので、ここではFortran90のMODULE文を使って渡します。⑧でPARAという名のMODULE文を作成し、その中でNPROCSとMYRANKを宣言します。さらにINCLUDE 'mpif.h'や、1対1受信サブルーチンで使用する配列ISTATUSなども、並列化する各サブルーチンで毎回宣言するのは面倒なので、MODULE文の中で宣言します。

そして、MODULE文内の変数を参照/更新したいサブルーチンは、MAIN文またはSUBROUTINE文の直後(宣言文より前)に、⑨と⑪に示すようにUSE文を指定するだけで、参照/更新が可能となります。なお、MODULE文は使い方に少しクセがあるので、4-5-7-1節「ブロック分割による配列の縮小方法3」の注意を厳守して下さい。

並列プログラムを1プロセスで実行すると、誤動作する可能性がありますので(4-8-2節の「プログラムの並列化方法のミス(1)」参照)、⑭でチェックを行い、1以下ならばエラーで終了します。

● プログラムが終了する前にMPI\_FINALIZEを一度だけ実行します。本例では⑬と⑯でプログラムが終了するので、⑯と⑰を追加します。3-2節で説明したように、MPI\_FINALIZEを指定せずに終了すると、マシンの環境によっては誤動作することがありますので、プログラムが終了する箇所が複数ある場合、全ての箇所ですべて一度MPI\_FINALIZEを実行してから終了するようにして下さい。

● 次に出力部分を修正します。②と⑤の出力部分(OPEN文やCLOSE文がある場合はそれも同様に)は、リンク⑱だけが出力するように、⑮と⑲のIF文を追加します(入出力部分の修正方法の詳細は4-4節を参照)。

④はプログラム内でエラーをチェックしているIF文です。最終的に並列化したプログラムで、このエラーを発見するのが全プロセスの場合⑳の修正になります。一方、このエラーを発見するのが特定のプロセスのみの場合、㉑だと、エラーを発見したプロセスのみがSTOP文で終了し、他のプロセスはIF文が偽になるので終了しません。特定のプロセスがエラーを発見したとき、全てのプロセスを終了させたい場合は、図4-8-1(3)の㉒に示すように、CALL MPI\_ABORT(付録参照)を実行して下さい。また特定のプロセスがリンク⑳以外の場合、㉓のIF文がついているとPRINT文が実行されないのので、㉔に示すようにIF文をはずします。また、どのプロセスがこのエラーを発見したのかを知りたい場合は、㉕の下線部を追加します。

これでとりあえず、プログラムが並列に実行できるようになりました。ただしこの時点では本来の意味の並列化をした部分はなく、プログラムを並列に実行すると、全プロセスが単体ジョブと全く同じ動作をし、ランク⑱のプロセスのみが出力を行います。

元の単体プログラムと上記の並列プログラムで、標準出力とその他の出力ファイルの内容が一致しているのを確認した後、本来の意味の並列化を行います。以下で概要を説明し、次の項で並列化を効率的に行うための方法を紹介します。

● ㉖(4-5-4節参照)で、⑦のループ反復のうち自己プロセスが担当する反復の下限(ISTA)と上限(IEND)を求め、それを㉗に追加し、㉘のようにループ反復を修正します。

なお、1つのプロセスが担当する要素数(本例ではIEND-ISTA+1)が少なすぎると、並列プログラムによっては誤動作しますので(4-8-2節の「プログラムの並列化方法のミス(2)」参照)、㉙で正常に動作する最低限の要素数(例えば2)以上かどうかをチェックし、それより少なければエラーで終了します(しきい値はプログラムのロジックに依存しますので、ロジックに応じて適切な値を設定して下さい)。

各プロセスが計算した結果を最後にランク⑱に集めるため、㉚の通信を行います。

```

PROGRAM MAIN
:
N = 1000
CALL SUB(N)
WRITE(10) X
STOP
END
SUBROUTINE SUB(N)
:
IF (エラーを発見したら)
PRINT *, '===ERROR==='
STOP
ENDIF
DO I=1,N
X(I) = A(I) + B(I)
ENDDO
:

```

図4-8-1(1)

```

MODULE PARA
INCLUDE 'mpif.h'
INTEGER ISTATUS(MPI_STATUS_SIZE)
INTEGER NPROCS, MYRANK, IOSTA, IEND
END
PROGRAM MAIN
USE PARA
:
N = 1000
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
IF (NPROCS<=1) THEN
IF (MYRANK==0)
& PRINT *, '===PARALLEL ERROR1===',
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
CALL PARA_RANGE
& (1, N, NPROCS, MYRANK, IOSTA, IEND)
IF (IEND-IOSTA+1<2) THEN
PRINT *, '===PARALLEL ERROR2===',
'MYRANK=', MYRANK
CALL MPI_ABORT
& (MPI_COMM_WORLD, 9, IERR)
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
CALL SUB(N)
各プロセスが計算した配列Xの結果を
ランク0に集める通信を行う。
IF (MYRANK==0) WRITE(10) X
CALL MPI_FINALIZE(IERR)
STOP
END
SUBROUTINE SUB(N)
USE PARA
:
IF (エラーを発見したら)
IF (MYRANK==0) PRINT *, '===ERROR===',
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
DO I=IOSTA, IEND
X(I) = B(I) + C(I)
ENDDO
:

```

図4-8-1(2) 全プロセスがエラーを発見する場合

```

:
IF (エラーを発見したら)
PRINT *, '===ERROR===', 'MYRANK=', MYRANK
CALL MPI_ABORT(MPI_COMM_WORLD, 9, IERR)
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
:

```

図4-8-1(3) 特定のプロセスだけがエラーを発見する場合

(6) 並列化

並列に動作させるための最低限の修正を行ったプログラムに対し、前述のように本当の意味の並列化を行います。並列化というのは通常のプログラム作成と感覚が違っており、慣れないうちは答えが合う状態に持っていくのがなかなか大変です。以下で、効率的に並列化するための方法をいくつか紹介します。

●【プロトタイプの作成】

大規模なプログラムや、並列化によって複雑な動作をするプログラムの場合、いきなり並列化するのではなく、実際のプログラムのうち並列化に関する本質的な特徴のみを抽出した100行以下の簡単な単体プログラムを作成し、それを並列化して、単体版と答えが同じになることを確認します。この過程で、本物のプログラムを並列化するために必要なノウハウ(例えば並列化に伴いプログラムのどこどどのような通信を行えば良いか、そのためにはどのようなようなテール類が必要か、など)を蓄積することができます。それに基づいて実際のプログラムを並列化すると、効率的に並列化することができます。

例えば図4-8-2(1)が10000行のプログラムで、そのうち計算時間のかかっている部分が①の2重ループだします。これに対し、以下の点を考慮して簡単化したプログラム(単体版)を図4-8-2(2)に示します。

- ①のループは1回計算するだけでテストできますので、タイムステップ・ループは削除しました。
- ①のループ以外の計算部分を削除しました。
- ①のループの2次元配列Pを1次元にし、サイズも小さくしました。
- ①のループで並列化に関する本質的な特徴は、 $P(I-1, J), P(I, J), P(I+1, J)$ を使用して $P(I, J)$ を計算していることなので、③ではその部分のみを抽出し、関係のない配列や変数(A, Bなど)は全て除去しました。
- 本物のデータを使用するのではなく、②に示すように人工的で簡単なデータを設定しました。これによってデバッグがしやすくなります。

```

:
入力データの読み込みと初期設定
DO ITIME=1, ITMAX (タイムステップループ)
:
EPS = 0.0
DO J=2, 999
DO I=2, 999
S = A(I, J)*P(I+1, J)+P(I, J+1))
+B(I, J)*(P(I+1, J+1)-P(I+1, J-1)
-P(I-1, J+1)+P(I-1, J-1))
+C(I, J)*(P(I-1, J)+P(I, J-1))
S = (S*D(I, J)-P(I, J))*E(I, J)
EPS = EPS + S*S
P(I, J) = P(I, J) + OMEGA*S
ENDDO
ENDDO
WRITE(10) P
END

```

図4-8-2(1) 10000行のプログラム

```

:
DO I=1, 10
P(I) = I
ENDDO
DO I=2, 9
P(I) = P(I-1)+P(I)+P(I+1)
ENDDO
DO I=1, 10
PRINT *, I, P(I)
ENDDO
END

```

図4-8-2(2) 100行以下のプログラム

●【データの動きが分かる図の作成】

1次元SOR法(5-6節参照)の例で説明します。

図4-8-3(1)を実行したときの配列Pの状態の推移を図4-8-3(2)に示します。○が計算前、□が計算後の要素、「境」は境界条件、↘と↙は要素の依存関係を表します。II=0のとき②,④,⋯を計算し、II=1のとき③,⑤,⋯を計算します。

これを並列化した図4-8-3(3)(4)では、まず↘の通信を行ってから②,④,⋯を計算し、次に↙の通信を行ってから③,⑤,⋯を計算します。

この例のように、並列化したときのデータの動きが複雑な場合、図4-8-3(2)(4)のように、データの動きが分かる図(要素数は必要最小限とし、プロセス数は必要最小限の3)を作成すると、並列化するときの参考になります。

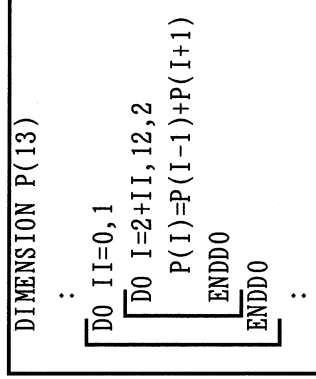


図4-8-3(1)

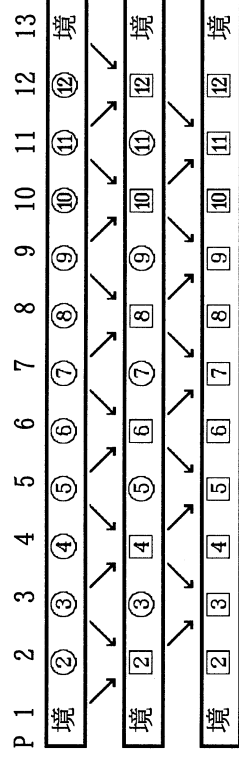


図4-8-3(2)

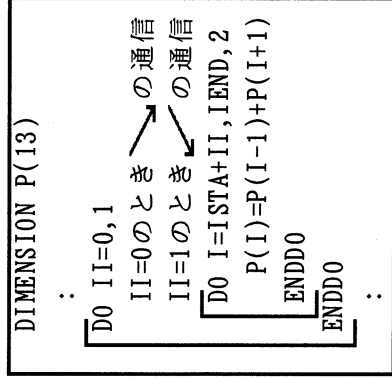


図4-8-3(3)

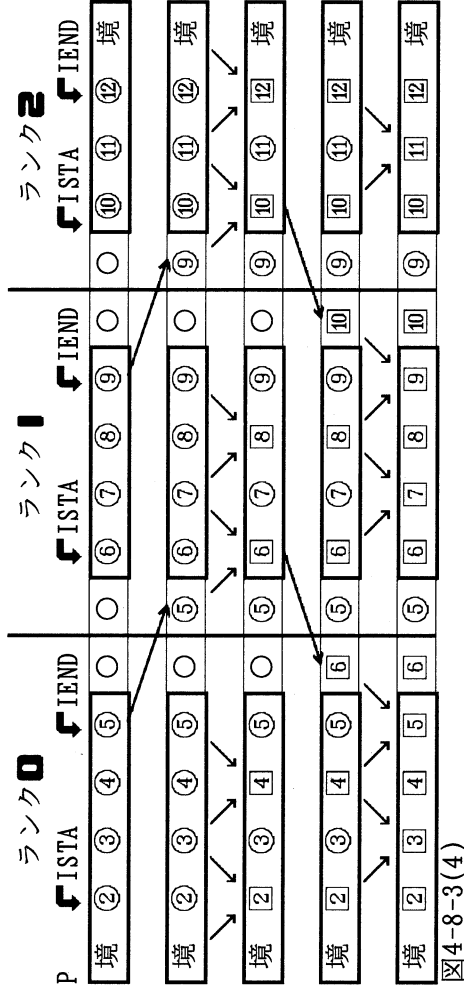


図4-8-3(4)

## ●【結果を確認しながら少しずつ並列化】

並列化に伴う修正箇所が多い場合(例えば差分法で現れる4-3節の「計算パターン2」)、並列化に伴う修正を全て行なってからデバッグするという方法では、バグがプログラムのあちこちに混入してしまい、バグを探すが『太平洋に落ちた針を探す』ように非常に難しくなります。

このような場合、少し並列化しては結果が正しいことを確認し、『だましました』少しずつ並列化の範囲を広げていくと、効率よく機械的に並列化することができます。この方法を図4-8-4(1)の例で説明します。

① 図4-8-4(1)のタイムステップ・ループ内に3つの(実際にはもっと多くの)D0ループがあり、これを全て並列化するとします。まず図4-8-4(2)のように1つ目のループを(本当の意味の)並列化します(ループ反復をISTA, IENDに修正)。それに伴って図4-8-5(1)のようなシフト通信が必要であれば、それを行うサブルーチンSHIFT(4-6-2節参照)をコールします。なお、図中で(本当の意味の)並列化した部分を『並』で表します。

1つ目のループを並列化したため、1つ目のループの終了後、各プロセスは、このループ内で値が設定された配列Aのうち、自分が計算した部分のデータしか持っていません。通常、配列Aは1つ目のループ以後の処理で参照/更新が行われるので、ここまでの修正を行っただけでは、並列に実行しても単体版と同じ計算結果にはなりません(4-3節の「計算パターン1」と同様)。

従って、この段階で正しい計算結果を得るためには、1つ目のループを終了した直後に、図4-8-5(2)に示すように、各プロセスは、配列Aのうち自分が計算した部分のデータを、他の全プロセスに送信する必要があります。そこで、これを行うサブルーチンDEBUG(4-6-3-2節参照)をコールします。これで各プロセスの配列Aには全てのデータが入り、単体プログラムと同じ状態になるため、並列に実行すると(バグがなければ)正しい計算結果が得られます。もし結果がおかしい場合は、並列化した1つ目のループ内か、サブルーチンSHIFT、DEBUGのいずれかにバグがあることとなります。このようにバグの存在する範囲が狭い範囲に限定されるので、バグを探すのが容易になります。なお、サブルーチンDEBUGは通信量が多いので、この時点でパフォーマンスは(一般に)却って低下します。

② 次に図4-8-4(3)に示すように2つ目のループを並列化し、ループ内で値が設定された配列WとBに対し、デバッグ用のサブルーチンDEBUGコールします。このとき、図4-8-4(2)で挿入したCALL DEBUG(A)を『並』部分の先頭に移動します。このように、CALL DEBUGは常に『並』部分の先頭に置くようにします。そして並列に実行し、結果が正しくなるまでデバッグします。

③ 図4-8-4(4)に示すように3つ目のループを同様に並列化し、ループ内で値が設定された配列Xに対し、デバッグ用のサブルーチンDEBUGをコールします。

配列Wは2つ目のループで値が設定され、3つ目のループで参照され、それ以外の部分では使われていないとします。つまり配列Wは、2つ目と3つ目のループ間でデータを受け渡す作業配列としてのみ使用されるとします。この場合、もはやCALL DEBUG(W)は不要なので、『並』部分の先頭に移動せずに削除します。CALL DEBUG(A)とCALL DEBUG(B)を『並』部分の先頭に移動して並列に実行し、結果が正しくなるまでデバッグします。

④ このように、結果が正しいことを確認しながら少しずつ並列化の範囲を拡大していくと、ついには図4-8-4(4)のようにタイムステップ内が全て『並』となり、しかも正しい結果が維持されています。ただしこの時点では、CALL DEBUGの影響でパフォーマンスは(一般に)かなり低下しています。

⑤ プログラムの最後に、配列Xの結果を10番ファイルに書き出しています。そこで例えば図4-8-4(5)に示すように、図4-8-5(3)の通信を行うサブルーチンGATHER(4-6-3-1節を参照)をコールし、ランク■が代表してファイルに書き出すようにし(他の出力方法については4-4-2節参照)、結果が正しくなるまでデバッグします。

最後にデバッグ用のCALL DEBUGを全て取り除くと、一気にパフォーマンスは向上し、結果も正しいはず…なのですが、もしここで結果がおかしくなくなった場合は、CALL DEBUGを(全てではなく)一つずつ取り除いていき、どれを取り除くと結果がおかしくなるのかを調べて下さい。

```

:
DO タイムステップ
DO I=1,N
:
A(I) = ~
ENDDO
DO I=1,N
:
W(I) = ~
B(I) = ~
ENDDO
DO I=1,N
:
~ = W(I)
X(I) = ~
ENDDO
WRITE(10) X
:

```

図4-8-4(1)

```

:
DO タイムステップ
並 CALL SHIFT(A)
並 DO I=ISTA, IEND
:
A(I) = ~
ENDDO
並 CALL DEBUG(A)
DO I=1,N
:
W(I) = ~
B(I) = ~
ENDDO
DO I=1,N
:
~ = W(I)
X(I) = ~
ENDDO
WRITE(10) X
:

```

図4-8-4(2)

```

:
DO タイムステップ
並 CALL SHIFT(A)
並 DO I=ISTA, IEND
:
A(I) = ~
ENDDO
並 CALL SHIFT(B)
並 DO I=ISTA, IEND
:
W(I) = ~
B(I) = ~
ENDDO
CALL DEBUG(A)
CALL DEBUG(W)
CALL DEBUG(B)
DO I=1,N
:
~ = W(I)
X(I) = ~
ENDDO
WRITE(10) X
:

```

図4-8-4(3)

```

:
DO タイムステップ
並 CALL SHIFT(A)
並 DO I=ISTA, IEND
:
A(I) = ~
ENDDO
並 CALL SHIFT(B)
並 DO I=ISTA, IEND
:
W(I) = ~
B(I) = ~
ENDDO
並 CALL SHIFT(X)
並 DO I=ISTA, IEND
:
~ = W(I)
X(I) = ~
ENDDO
CALL DEBUG(A)
CALL DEBUG(B)
CALL DEBUG(X)
ENDDO
WRITE(10) X
:

```

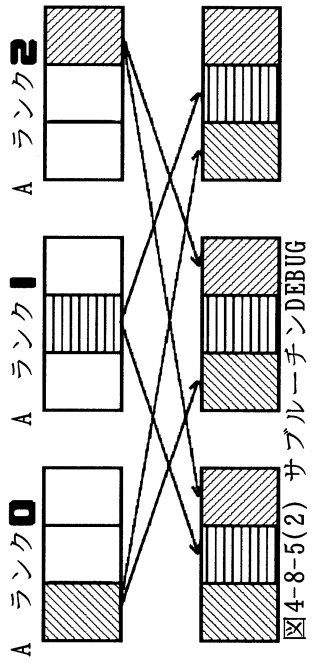
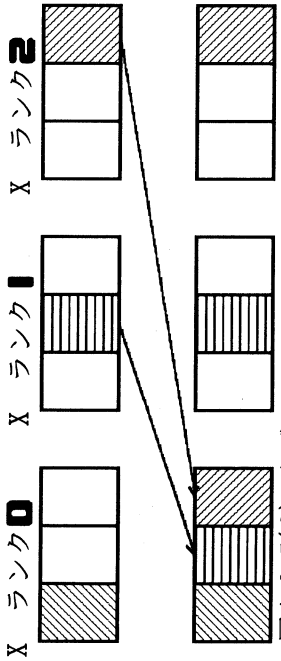
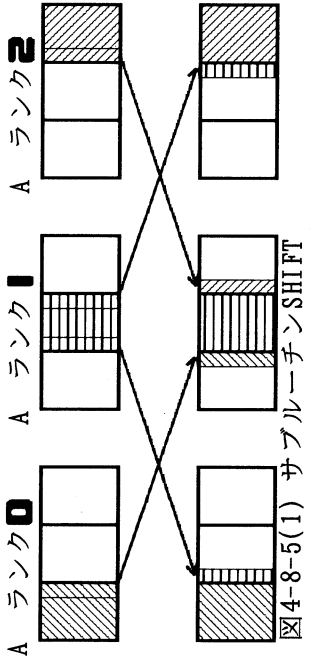
図4-8-4(4)

```

:
DO タイムステップ
並 CALL SHIFT(A)
並 DO I=ISTA, IEND
:
A(I) = ~
ENDDO
並 CALL SHIFT(B)
並 DO I=ISTA, IEND
:
W(I) = ~
B(I) = ~
ENDDO
並 CALL SHIFT(X)
並 DO I=ISTA, IEND
:
~ = W(I)
X(I) = ~
ENDDO
CALL DEBUG(A)
CALL DEBUG(B)
CALL DEBUG(X)
ENDDO
CALL GATHER(X)
IF (MYRANK==0)
WRITE(10) X
ENDIF
:

```

図4-8-4(5)





プログラムを並列化した場合の典型的なエラーとその対処法について説明します(本文ですでに説明済みの項目もありますが、念のため再掲します)。なお、下記の内容はマシン環境によっても異なります。

#### ■ システムのトラブル

並列プログラムは正しいはずなのに実行すると挙動がおかしい場合、システムのトラブルの可能性がありますが、3-2節の図3-2-1(1)のような簡単なプログラムを実行し、システムが正常に動作しているかどうかをまず確認して下さい。

#### ■ 実行方法のミス

各ノードのローカルディスクに置かれたロードモジュールa.outを使用して並列ジョブを実行する環境では、並列ジョブの実行前にあらかじめa.outを各ノードのローカルディスクにコピーします。例えば並列プログラムの誤りを修正し、再コンパイル・リンクして実行したところ、結果が修正前と同じでおかしいなど思い、原因を調べたところa.outをコピーし忘れていたという状況が発生しますので、コピーのし忘れには注意して下さい。

#### ■ 結果の合法的な相違(誤差など)

- 並列化に伴う修正によって、その付近のロジックに対するコンパイラの最適化の度合いが変わり、計算結果が変わることがあります。この場合、試みにコンパイルオプションを指定せずに単体版プログラムと並列版プログラムを実行し、もし結果が同一になれば、(コンパイラのバグでない限り)最適化の度合いが変わったことが原因であることとなります。
- 合計や内積を求めるD0ループを並列化した場合、単体で実行したのと計算順序が変わるので、誤差によって合計の結果も若干変わります(3-3-6節参照)。敏感なプログラムの場合、これがその後の計算結果に影響することがあります。これは合法的な誤差ですが、プログラムが単精度の場合には倍精度にすれば回避できる場合もあります。
- 乱数を使用したプログラムを並列化した場合は、発生した乱数系列が単体版のプログラムと(通常)変わるので、結果も変わります。

#### ■ MPI自体のバグ

MPIのルーチン自体にバグがある可能性があると言えないとは言えません(例えば集団通信ルーチンの内部でデッドロックしてしまうなど)。MPIのルーチンがおかしいと思われる場合、もし他のMPIを使用するマシン環境があればそこで実行し、同様の現象が起こらなければ、最初のマシン環境でのMPIのバグだと思われま

す。その場合、電算部門やメーカーの担当者にお問い合わせますが、その際プログラムや入出力データを、なるべく以下のように簡単化してから問い合わせた方がよいでしょう。その目的は、担当者がバグの解析をしやすく回答が速く得られることと、簡単化する過程で、例えばプログラムのこの部分を削除したらエラーが起きなくなっただなど、エラーの手がかりが見つかることがあるためです。

- プログラムを、問題が発生する最低限の行数(出来れば数10行)に縮小する。
- 入出力ファイルの数やデータ量を最小限(出来ればゼロ)にする。
- 問題が発生するまでの実行時間を最短(出来れば数秒)にする。

#### ■ 元の単体プログラムに含まれている潜在的なバグ

並列化する前の単体プログラムに含まれている潜在的なバグ(例えば配列の範囲を越えて代入を行ったため記憶域を壊していたなど)が、並列化したことによって計算結果の相違として表面に現れる場合があります。その場合、単体版または並列版のプログラムにデバッグ用のコンパイル・リンクオプション(例えばモジュール間の整合性をチェックするオプション)を付けて実行すると、その箇所を発見できる場合があります。あるいは上の項目で説明した、プログラムを簡単化する過程でその箇所を発見できる場合もあります。

## ■ マシン環境によるMPIの制限

- (4-4-1節) 標準入力(装置番号5番)から読み込みを行うと、マシン環境によってはエラーになる場合があるのでお勧めできません。
- (3-3-4節、3-4-5節) 1対1通信および集団通信ルーチンで一度に送信できる最大のバイト数はマシン環境によって異なり、無制限の場合と、制限のある場合があります。制限のある場合に、制限を越えたバイト数を送ってもエラーにならず、不完全なメッセージが送られてしまうマシン環境もあります。
- (3-4-5節) 1対1通信ルーチンで仕掛かり中のメッセージがどんどん累積していった場合、マシン環境によっては例えば受信バッファがパンクして異常終了するなどの現象が発生します。
- MPIの使い方のミス
- MPIの使い方自体に不明な点がある場合は、参考文献[15]に記載したアドレスにメール(ただし英語)すると、回答が得られる場合があります。
- (3-2節) 「MPI\_」で始まる変数名、関数名、サブルーチン名をユーザープログラムで使用すると、MPIが使用する変数と偶然一致して誤動作する危険性がありますので、使用しないで下さい。
- (3-2節) MPI\_FINALIZEをコールせずに終了した場合、マシン環境によっては誤動作する(実行中の他のプロセスを強制的に終了させてしまう、あるいは実行するたびに計算結果が異なる再現性のないエラーとなる)ことがあります。プログラムが終了する全ての箇所(STOP文またはEND文)で、必ず1度だけMPI\_FINALIZEをコールしてから終了するようにして下さい。
- (3-3-4節) MPIのサブルーチンの引数が少なかったり指定がおかしくても、一般にコンパイル・リンクではエラーになりません。この場合、実行時にSegmentation faultなどのエラーで異常終了しますが、エラーメッセージを見ても原因が分からない場合がありますので、異常が起きた場合にはMPIのサブルーチンの引数の数や属性に誤りがないかどうかをチェックして下さい。特に多いのは以下のエラーです。
  - CALL MPI\_xxx(～, IERR)の整数IERRを指定し忘れた場合。
  - MPI\_RECV, MPI\_WAITで使用する引数statusを、INTEGER ISTATUS(MPI\_STATUS\_SIZE)で宣言し忘れた場合。
  - INCLUDE 'mpif.h'を指定し忘れた場合。
- (3-3-4節) MPIの仕様書(参考文献[15])では、ユーザーがMPIの使い方を間違えた場合、実行時にMPIが行うエラー処理について一般に規定していません。このため、例えばMPI\_BCASTなどの集団通信ルーチンで、送信バッファの大きさより受信バッファの大きさを誤って小さくしてしまったような場合でも、受信時にエラーにならず、エラーのリターンコード(ierror)にも正常終了の「MPI\_SUCCESS」が返るマシン環境もあります。
- (3-4-2節) 非ブロッキング通信のMPI\_ISENDまたはMPI\_RECVを指定してMPI\_WAITを指定し忘れた場合、タイミングによる(再現性のない)エラーが発生します。またMPI\_ISENDまたはMPI\_RECVとMPI\_WAITの引数requestがスเปルミスで異なっている場合も、当然ながらMPI\_WAITは働かず、タイミングによる(再現性のない)エラーが発生しますので、スเปルミスには十分注意して下さい(一方のスペルをカット・アンド・ペーストで他方にコピーするのが安全です)。
- (3-4-4節) 双方向通信を行う際、1対1通信ルーチンを実行する順序がおかしいと、デッドロックが発生してジョブが止まってしまいます。デッドロックになっているかどうかは、通信の前後にPRINT文を入れてランク番号を書き出し、各プロセスがその通信を終了しているかどうかで調べることができます。
- (3-4-5節、3-3-4節) 1対1通信ルーチンでは、受信バッファで指定するバイト数(データ個数×データ型のバイト数)は、送信バッファで指定するバイト数と等しいか、または大きくても構いません。一方集団通信ルーチンでは、送受信バッファで指定するバイト数は一致していなければなりません。

■ プログラムの並列化方法のミス

(1) 並列プログラムを1プロセスで実行すると、誤動作する可能性がありますので(3-4-5節の「1対1通信を1プロセスで実行した場合」、4-6-2節の「境界条件の処理」参照)、1プロセスで実行するときは、並列版ではなく単体版のプログラムを使用することを強くお勧めします。並列プログラムの開始時に、全プロセス数が1以下の場合はエラーで終了させるようにして下さい(4-8-1節の「(5) 並列に動作させる最低限の修正」参照)。

(2) 図4-8-3(2)のプログラムでは、図4-8-3(1)に示すように、配列Aの3つの要素から配列Bの要素を計算します。これを並列化した図4-8-3(3)では、各プロセスはブロック分割(4-5-4節参照)で決定したISTAからIENDまでの要素を担当します。これを2プロセスで並列に実行した場合、図4-8-3(4)に示すように、ランク0の右端の2要素をランク1のプロセスに通信(\\の部分)してから計算に入ります。

3プロセスで並列に実行した場合、図4-8-3(5)に示すように、ランク0とランク1のプロセスの右端の2要素を右隣のプロセスに通信します。ところがランク1のX(値が未定義)の要素がランク2に送られ、ランク2はXの要素を使用して計算を行うため、誤った計算結果になってしまいます。

これはランク1のプロセスが担当する要素数が1つしかないのが原因です。つまりこの例では、各プロセスが担当する要素数が最低2以上でないと、正常に動作しません。

このように、1つのプロセスが担当する要素数が少なすぎると、並列プログラムによっては誤動作します。その並列プログラムで、1つのプロセスが担当する要素数が最低いくつあれば正常に動作するかを把握し、プログラムの実行を開始して各プロセスが担当する要素数が決定した時点で、正常に動作する最低限の要素数以上かどうかをチェックするようにして下さい(4-8-1節の「(5) 並列に動作させる最低限の修正」参照)。

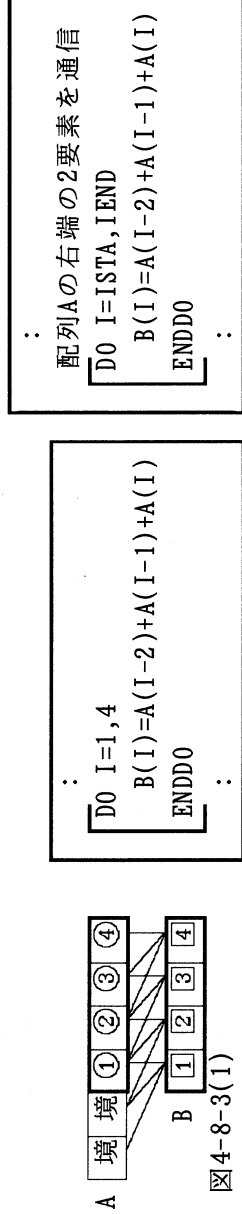


図4-8-3(1)

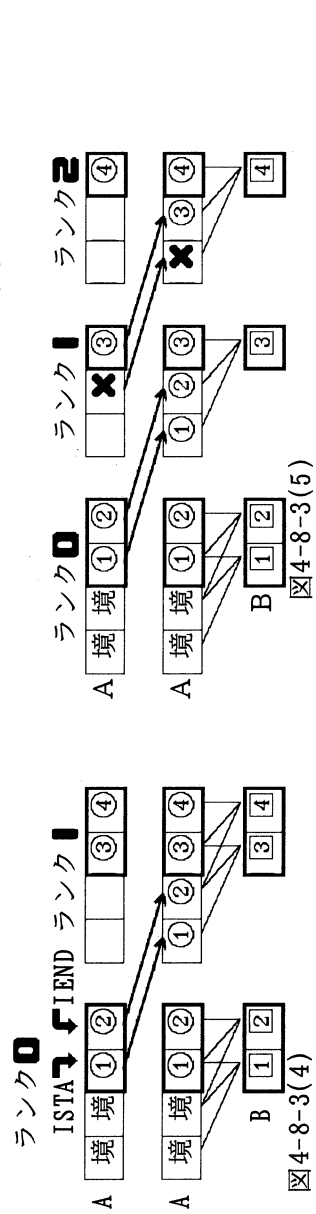


図4-8-3(2)

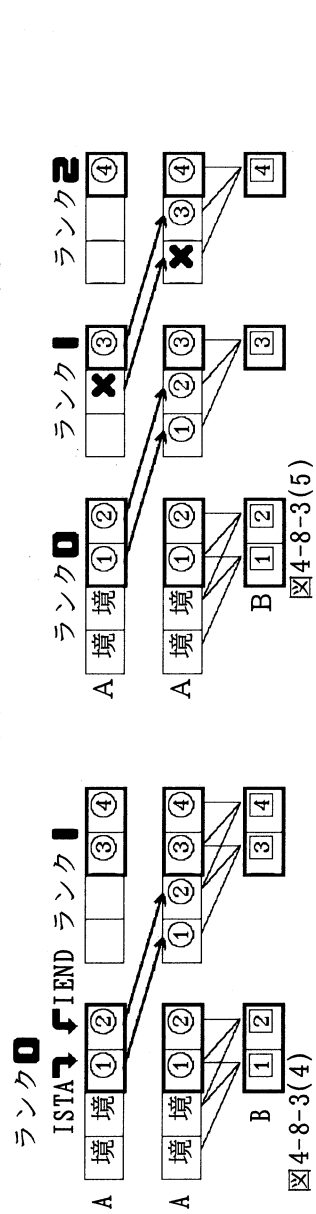


図4-8-3(3)

図4-8-3(4)

(3) 4-5-4節で説明したブロック分割がうまくいっていないため、各プロセスが担当するループ反復の下限(ISTA)と上限(IEND)の値がおかしい場合があります。全プロセスのISTAとIENDの値をPRINT文などで書き出し、確認して下さい。

(4) 並列性のないループ(回帰参照のあるループ)をうまくいって並列化した場合、当然ながら結果がおかしくなります。単体版プログラムでそのループを逆に反復させてみて、元の計算結果と変わるようであれば、一般にそのループに並列性はありません(ただし例外もあります)。

## 4-8-3 パフォーマンスの測定と評価

ここまでで一応並列化が完了しました。実は、並列化には大きく以下の2つの関門があり、まだ(1)の関門を突破したに過ぎないのです。

- (1) 並列化して結果が合うこと
- (2) 良好なパフォーマンスが得られること

本節では、(2)に関係のある、パフォーマンスの測定および評価方法について説明します。

### ■ 並列ジョブの評価方法

単体ジョブの場合、主にCPU時間で測定を行います。並列ジョブの場合は以下のように通信時間と通信待ち時間が加わるため、経過時間で測定を行います。

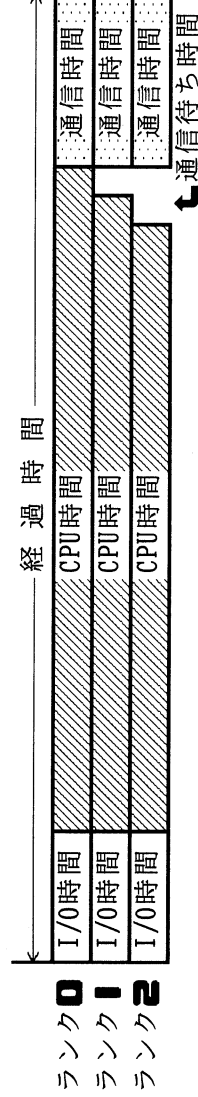


図4-8-5

並列ジョブの効果を評価する方法として、以下に示す速度向上率がよく使用されます。例えば4台で実行した場合、完全に並列化できた場合には速度向上率が4倍となりますが、通常は通信などのオーバーヘッドで効率が低下し、4倍以下になります。

$$\text{速度向上率} = (\text{1台での経過時間}) / (\text{n台での経過時間})$$

速度向上率に関する注意点を以下に示します。

- 分子の(1台での経過時間)は、並列版のプログラムをプロセス数1として実行した場合の経過時間ではなく、単体のプログラムで測定した経過時間にして下さい。また経過時間は、他のジョブが流れていない占有環境で測定して下さい。
- ロードモジュールa.outや入力データをローカルディスクに置いた場合と共有ディスクに置いた場合で経過時間が異なります。特にノード数が多い場合、共有ディスクを使用すると多くのノードが同じファイルアクセスして競合するため、最初のノードにロードモジュールがロードされた時刻と最後のノードにロードされた時刻が大きく離れてしまう場合があります。
- 一般的なプログラムでは処理は以下のような流れになります。

- (1) 入力データの読み込みや初期設定
- (2) 計算部分
- (3) 計算結果の出力

通常は(2)の部分が経過時間のほとんどを占め、(1)と(3)の割合がわずかです。ところが、測定時間を短縮するために通常よりも短いジョブで測定した場合、(1)や(3)の相対的な割合が高くなってしまいます。さらに、(1)や(3)は通常並列化しないのでますます割合が高くなってしまいます。このような場合、(2)の部分だけをプログラム内でタイマーを使用して測定した方がよいでしょう。タイマーを使用した測定方法については後述します。

● 速度向上率はあくまで相対的な尺度であり、絶対的な評価基準はあくまで経過時間そのものであることに注意して下さい。

例えば図4-8-6(1)に示すように、ある単体プログラムのCPU時間が400分かかっており、計算部分全体をロードバランスが均等になるように並列化してきました。これを4ノードで実行したところ、CPU時間は100分(着色部分)になります。また並列化に伴い、通信時間が20分(白い部分)かかったとします。この場合、速度向上率は3.3倍と、まずまずの値となります。

このプログラムを単体チューニングし、図4-8-6(2)に示すように、単体プログラムのCPU時間が1/4の100分になったとします。これを並列に実行すると、通信時間は単体チューニング前と変わらないので、経過時間は45分になり、速度向上率は2.2倍にしかなりません。

このように、単体チューニングをすると、経過時間自体は短くなりますが、速度向上率は低下します。

一方、ある並列プログラムをある並列計算機で実行したところ、図4-8-6(2)のように速度向上率が2.2倍だったとします。同じ並列プログラムを、通信性能が同じでCPU性能が4倍遅い並列計算機で実行すると、図4-8-6(1)のように速度向上率は3.3倍になります。

このように、同じ並列プログラムの場合、通信性能の低い並列計算機の方が、速度向上率は高くなります。また同じ理屈で、同じ並列プログラムの場合、CPU性能が同じであれば、通信性能の高い並列計算機の方が、速度向上率は高くなります。

まとめると、同じ並列プログラムでも、並列計算機のCPU性能と通信性能の比率によって、速度向上率は変動します。

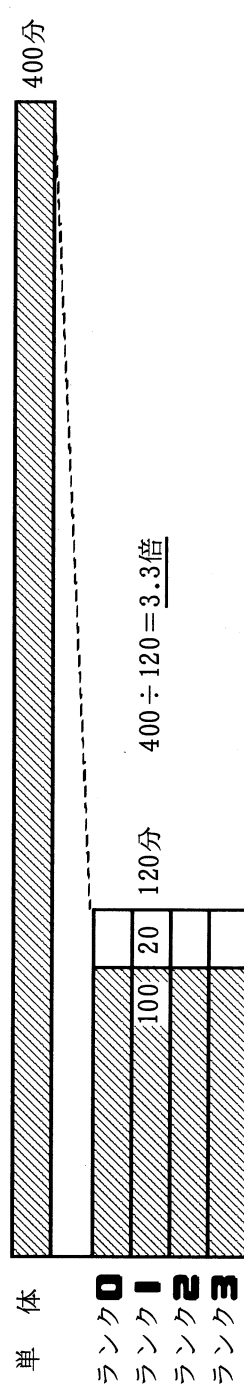


図4-8-6(1) 単体チューニング前/CPU性能の低い並列計算機

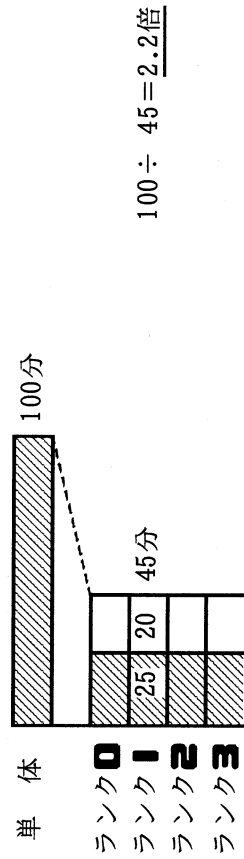


図4-8-6(2) 単体チューニング後/CPU性能の高い並列計算機

■ 経過時間の測定方法

図4-8-7の①のようにtimeコマンド(またはtimexコマンド)を指定してジョブを実行すると、②のような結果が表示されます(表示内容はシェルによって異なります)。realはジョブ全体の経過時間(秒)を示します。userは、単体ジョブの場合はCPU時間ですが、並列ジョブの場合、マシン環境によって表示が異なり、図のようにほとんどゼロの値が表示される場合と、全プロセスのCPU時間の合計が表示される場合があります。図で、マニュアルなどで確認して下さい。

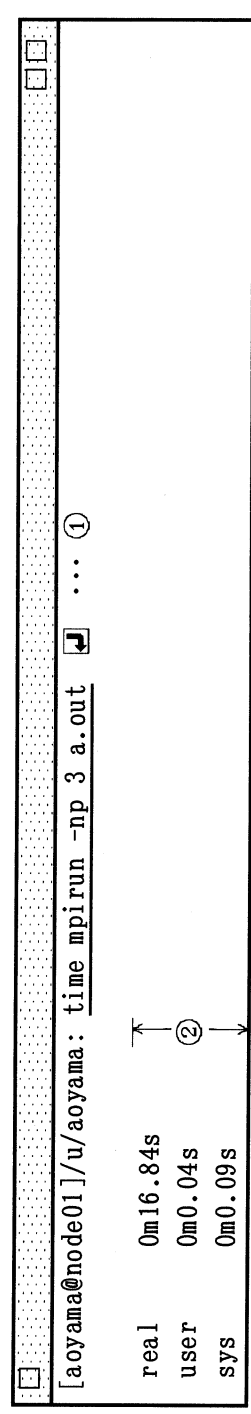


図4-8-7

## ■ プログラム内の特定の部分の経過時間の測定

並列化したプログラム内の特定の部分の経過時間を測定したい場合、MPIのファンクションMPI\_WTIME()を使用します(「付録」参照)。図4-8-8(1)に示すように、「測定対象部分」の前後で②と③を実行すると、変数ELP1とELP2(変数名は任意)に(内部的な)時刻が戻ります。④で両者の差をとれば、測定部分の経過時間(単位は秒)が求まります。なお、MPI\_WTIME()は倍精度なので、①の指定を必ず行って下さい。図4-8-8(1)を3プロセスで実行した場合のタイムチャートを図4-8-8(2)に示します。

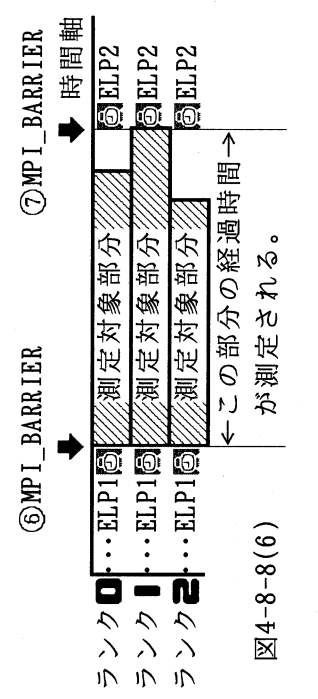
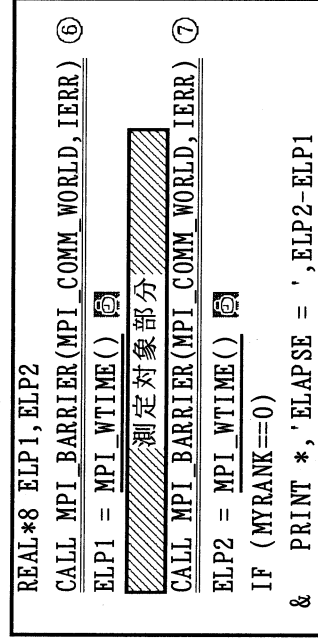
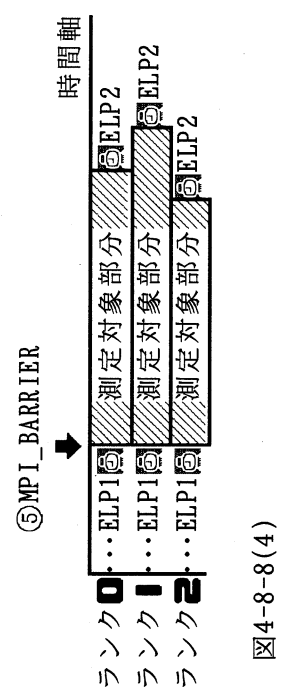
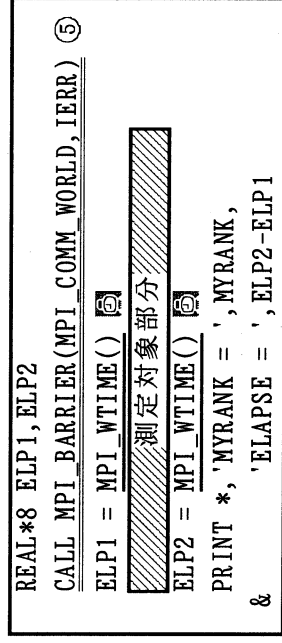
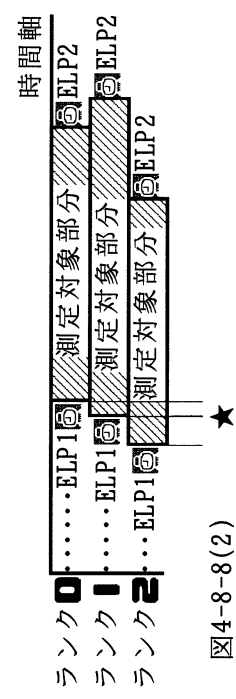
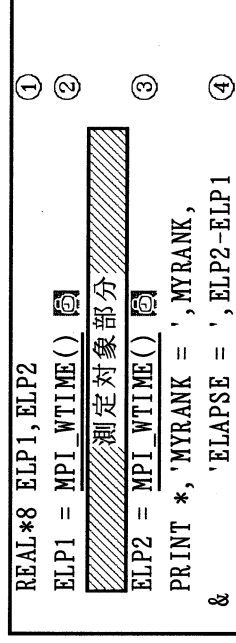
ところで、並列プログラムのロードモジュールa.outが共有ディスクにある場合、a.outのロードが完了してプログラムが開始する時刻がプロセスによって若干異なります。また各プロセスが同一ファイルから入力データを読み込む場合も、読み込みを完了する時刻がプロセスによって異なります。このように、プログラムが開始した直後は一般に各プロセスの進み具合が若干異なっており、その影響で、図4-8-8(2)の★に示すように、各プロセスが「測定対象部分」に到達する時刻が若干異なることがあります。

この影響をなくすためには、図4-8-8(3)の⑤に示すように集団通信サブルーチンMPI\_BARRIER(「付録」参照)を使用します。すると⑤の時点で同期が取られるため、タイムチャートは図4-8-8(4)のようになり、「測定対象部分」の開始時点でのバラつきはなくなります。

例えば、DOLRUPをプロセス分割で並列化し、各プロセスが担当している計算量のロードバランスが均等になっているかどうかを調べたい場合にこの方法を用います。

一方、図4-8-8(5)の⑥、⑦に示すように、「測定対象部分」の前後にMPI\_BARRIERを使用すると、タイムチャートは図4-8-8(6)のようになり、「測定対象部分」の最長の経過時間が測定されます。

例えば、プログラムの主要計算部分(プログラムの最初と最後に1度だけ行う入出力や初期設定を除いた部分)のみの経過時間を調べたい場合にこの方法を用います。



■ 並列化したプログラムのパフォーマンスが悪い場合の調査方法

並列化したプログラムのパフォーマンスが悪い場合、タイマーチェーンを使用して原因を調べます。例えば図4-8-9のプログラムで、①は(本当の意味の)並列化をしていない部分、②は(本当の意味の)並列化をした部分、③は並列化に伴って必要となった通信部分だとします。①,②に示すように、各部分の前後に前述のタイマーチェーンMPI\_WTIME()を挿入して経過時間を測定し、それを③でEELP(1)~EELP(3)に累積します(①,②,③が複数箇所ある場合は、それぞれの場合、それぞれ箇所で測定を行います)。

①と③では、一番長かったプロセスの経過時間を測定したいので、④,⑤に示すように、1回目と2回目のMPI\_WTIME()の前にMPI\_BARRIERを挿入して同期を取ります。一方②では、各プロセスのロードバランスのばらつきを調べたいので、⑥に示すように、1回目のMPI\_WTIME()の前だけにMPI\_BARRIERを挿入します。

そして他のジョブが流れていない占有使用の状態、例えば1,2,4プロセスで測定を行います。なお、1プロセスでの測定は、並列版ではなく単体版のプログラムを使用します(ただし単体版ではMPI\_WTIME()が使用できないので、コンパイラが提供している経過時間測定ルーチンを使用します)。

並列ジョブのパフォーマンスが悪い場合、4-1-2節で説明したように、以下の3つの原因が考えられますので、測定結果からどのケースに該当するのかを調べ、その問題を(可能であれば)取り除きます。

- 図4-8-10(1)では、並列化率が低過ぎるため、並列化しても速くなりません。
- 図4-8-10(2)では、各プロセスのロードバランスが不均等なため、並列化しても速くなりません。
- 図4-8-10(3)では、通信に時間がかかっているため、並列化しても速くなりません。

```

:
REAL*8 ELP1, ELP2, EELP(3)
:
DO I=1,3
  EELP(1) = 0.0
ENDDO
DO ITIME=1, ITMAX  タイムステップ・ループ
  CALL MPI_BARRIER(MPI_COMM_WORLD, IERR) ④
  ELP1 = MPI_WTIME() ①
DO I=1,N
  :
  ENDDO
  CALL MPI_BARRIER(MPI_COMM_WORLD, IERR) ⑤
  ELP2 = MPI_WTIME() ②
  EELP(1) = EELP(1)+ELP2-ELP1 ③
  CALL MPI_BARRIER(MPI_COMM_WORLD, IERR) ⑥
  ELP1 = MPI_WTIME() ①
DO I=1,STA, IEND
  :
  ENDDO
  ELP2 = MPI_WTIME() ②
  EELP(2) = EELP(2)+ELP2-ELP1 ③
  CALL MPI_BARRIER(MPI_COMM_WORLD, IERR) ④
  ELP1 = MPI_WTIME() ①
  CALL MPI_BCAST(~) ③通信部分
  CALL MPI_BARRIER(MPI_COMM_WORLD, IERR) ⑤
  ELP2 = MPI_WTIME() ②
  EELP(3) = EELP(3)+ELP2-ELP1 ③
ENDDO
PRINT *, 'MYRANK=', MYRANK, EELP
:

```

図4-8-9

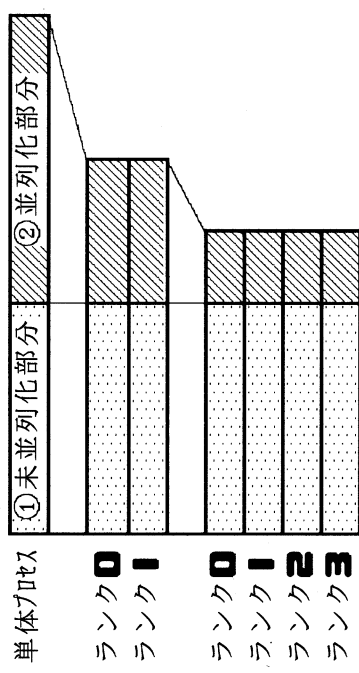


図4-8-10(1) 並列化率が低いケース

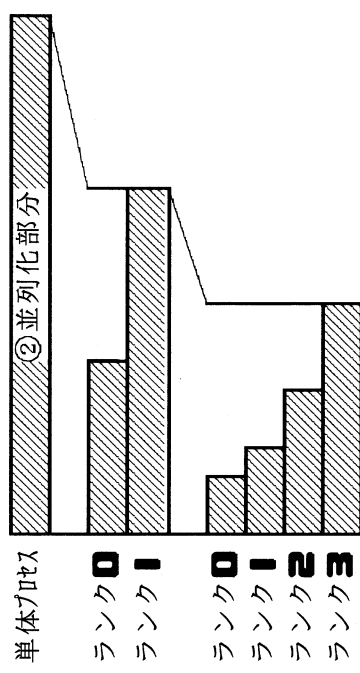


図4-8-10(2) ロードバランスが不均等なケース

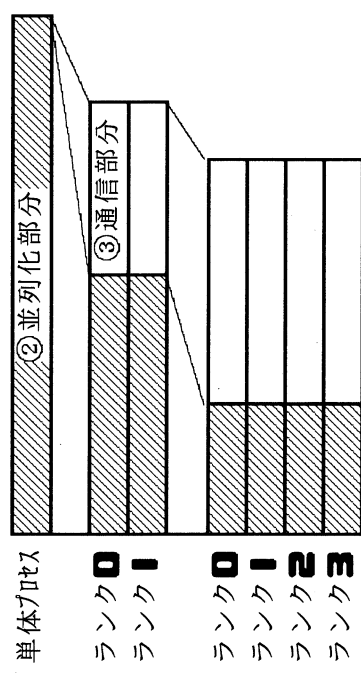


図4-8-10(3) 通信が多いケース

(このページは空白です。)



## 第5章

### 並列化の応用例

第4章までで、プログラムを並列化するための技法についての説明は終わりました。本章では、データの動きが複雑なプログラムの並列化例、ややマニアックな並列化例を紹介します。

## 5-1 差分法徹底攻略

キング・オブ科学技術計算の称号を差分法に与えても異議をとねえ人はいないでしょう。実際、過去に筆者が並列化したプログラムのうち、約半数が差分法のプログラムでした。差分法で行うシフトについては、4-6-2節で1次元差分法モードキを例に基本的な考え方を説明しました。本節では2次元差分法でシフトを行う際に必要となる各種のテクニクを徹底攻略します。

なお、シフトの通信方法には、4-6-2節の図4-6-3(1)~(3)に示す「通信方法1」~「通信方法3」の3種類があり、どの方法が最も速いかはマシン環境によって異なります。本節では便宜的に「通信方法1」を使用します。

また、3次元差分法は2次元差分法での方法を単に拡張するだけなので説明は省略します。

4-6-2節の1次元の差分法モードキのプログラム(図4-6-2(1))を2次元に拡張したのが図5-1-1(1)です。まず①で配列Aに適当な値を設定し(図5-1-1(2)の○と●)、②で配列Aの上下左右の値(●)から配列Bの値(■)を求めます。

```

PROGRAM MAIN
  IMPLICIT REAL*8(A-H,0-Z)
  PARAMETER (M=4, N=9)
  DIMENSION A(M,N), B(M,N)
  DO J = 1, N
    DO I = 1, M
      A(I,J) = I + 10.0*J
    ENDDO
  ENDDO
  DO J = 2, N-1
    DO I = 2, M-1
      B(I,J) = A(I-1,J) + A(I,J-1)
      &          + A(I,J+1) + A(I+1,J)
    ENDDO
  ENDDO
  END
  
```

図5-1-1(1)

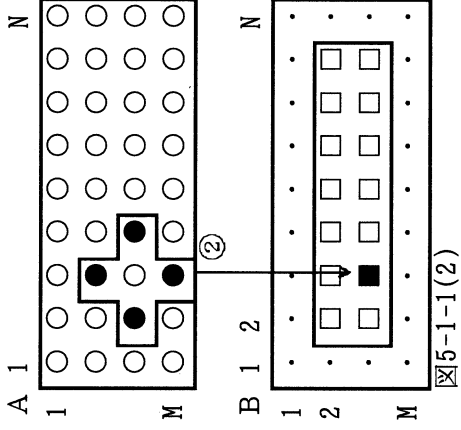
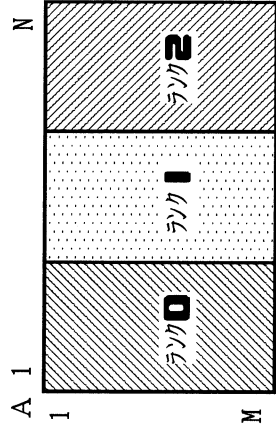


図5-1-1(2)

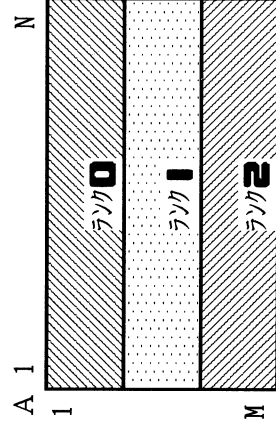
2次元の差分法プログラムを並列化する場合、配列A, Bの分割には以下の3つの方法があり、通信量などの点からどれを選ぶかを決めます(4-5-8節参照)。以下の各節では、それぞれの分割の具体的な方法について説明します。



2次元目でブロック分割

→ 5-1-1節で説明

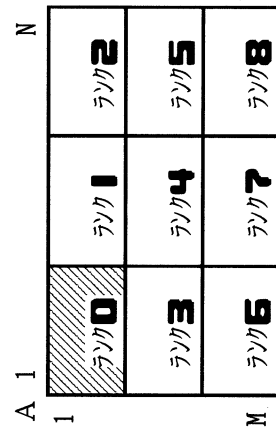
図5-1-2(1)



1次元目でブロック分割

→ 5-1-2節で説明

図5-1-2(2)



1,2次元目でブロック分割

→ 5-1-3節で説明

図5-1-2(3)

## 5-1-1 2次元目でブロック分割した場合

本節では、配列A, Bを図5-1-2(1)のように2次元目でブロック分割する場合のシフト方法を説明します。並列化したプログラムを図5-1-3(1)に、データの動きを図5-1-3(2)に示します(配列AとBのうち、各プロセスが実際に使用する部分を太線で示します)。図4-6-3(1)(2)と図5-1-3(1)(2)と図5-1-3(1)内の番号は対応しています。

図4-6-3(1)ではシフトするデータが1つの要素でしたが、図5-1-3(1)では1列全体をシフトするので、⑥～⑨で送受信するデータの要素数は『M』になります。それ以外は図4-6-3(1)と同じなので、これ以上の説明は省略します。

```

PROGRAM MAIN
IMPLICIT REAL*8(A-H,0-Z)
INCLUDE 'mpif.h'
PARAMETER (M=4,N=9)
DIMENSION A(M,N),B(M,N)
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD,NPROCS,IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD,MYRANK,IERR)
CALL PARA_RANGE
& (1,N,NPROCS,MYRANK,JSTA,JEND)
JSTA2 = JSTA
JEND1 = JEND
IF (MYRANK == 0 ) JSTA2 = 2
IF (MYRANK == NPROCS-1) JEND1 = N-1
IUP = MYRANK + 1
IDOWN = MYRANK - 1
IF (MYRANK==NPROCS-1) IUP =MPI_PROC_NULL
IF (MYRANK==0) IDOWN=MPI_PROC_NULL
DO J = JSTA, JEND
DO I = 1, M
A(I,J) = I + 10.0*J
ENDDO
ENDDO

```

```

CALL MPI_ISEND(A(1,JEND) ,M,
& MPI_REAL8,IUP ,1,
& MPI_COMM_WORLD,ISEND1,IERR)
CALL MPI_ISEND(A(1,JSTA) ,M,
& MPI_REAL8,IDOWN,1,
& MPI_COMM_WORLD,ISEND2,IERR)
CALL MPI_IRECV(A(1,JSTA-1),M,
& MPI_REAL8,IDOWN,1,
& MPI_COMM_WORLD,IRecv1,IERR)
CALL MPI_IRECV(A(1,JEND+1),M,
& MPI_REAL8,IUP ,1,
& MPI_COMM_WORLD,IRecv2,IERR)
CALL MPI_WAIT(ISEND1,ISTATUS,IERR)
CALL MPI_WAIT(ISEND2,ISTATUS,IERR)
CALL MPI_WAIT(IRecv1,ISTATUS,IERR)
CALL MPI_WAIT(IRecv2,ISTATUS,IERR)
DO J = JSTA2, JEND1
DO I = 2, M-1
B(I,J) = A(I-1,J) + A(I,J-1)
+ A(I,J+1) + A(I+1,J)
ENDDO
ENDDO
CALL MPI_FINALIZE(IERR)
END

```

図5-1-3(1)

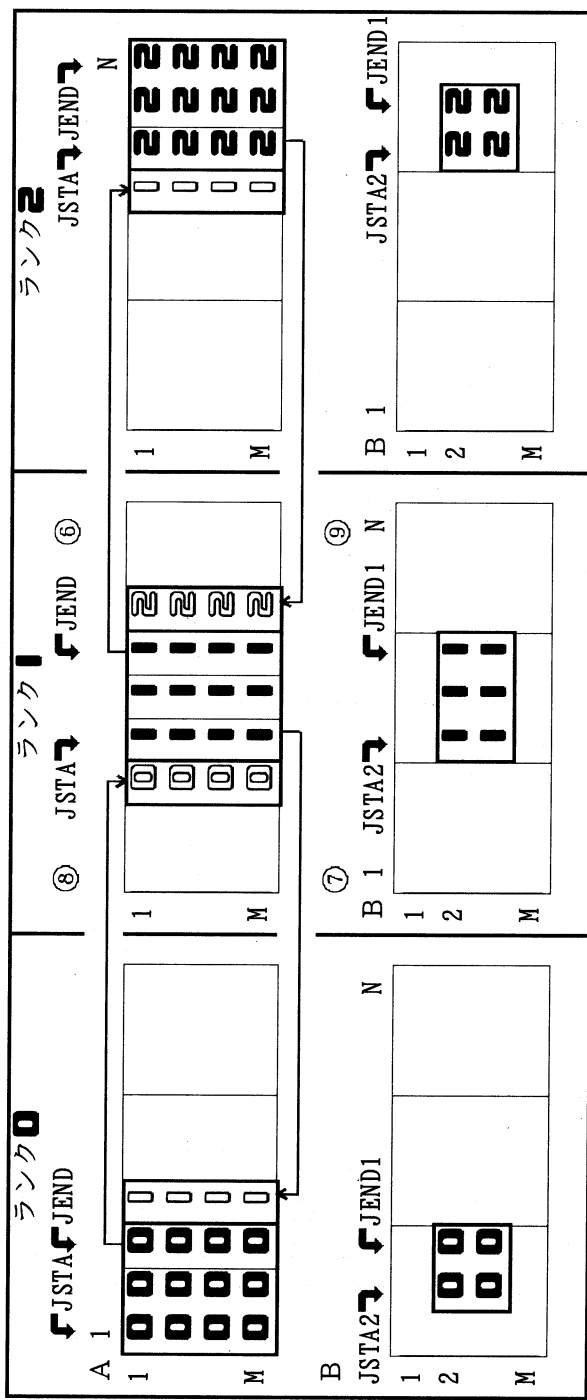


図5-1-3(2)

## 5-1-2 1次元目でブロック分割した場合

配列A, Bを図5-1-2(2)のように1次元目でブロック分割する場合のシフト方法を説明します(紙面の関係で配列A, Bの大きさは前節と異なり $9 \times 4$ とします)。並列化したプログラムを図5-1-4(1)に、データの動きを図5-1-4(2)に示します。図5-1-4(1)(2)内の番号はそれぞれ対応しています。

配列A, Bを1次元目でブロック分割した場合、通信するデータはプロセス間の境界の1行ですが、1行内の各要素はメモリー上でとびとびになっているので、簡単には通信できません。そこで派生データ型を使用して通信を行います。まず⑫で、配列A(と同じ大きさの配列)の任意の行の先頭から始まる長さNの1行を表す派生データ型ITYPEを、3-5-3-2節で説明した自作のサブルーチンPARA\_TYPE\_BLOCK2Aを使用して作成します(MPI-2が使用できるマシン環境の場合は、MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照))を使用できますが、派生データ型の起点が異なるので注意して下さい)。そして⑥~⑨の下線部で、派生データ型ITYPEを指定して実際の送受信を行います。派生データ型を使用する以外の部分は前節の例と同じなので、これ以上の説明は省略します。

```

PROGRAM MAIN
IMPLICIT REAL*8(A-H,0-Z)
INCLUDE 'mpif.h'
PARAMETER (M=9,N=4)
DIMENSION A(M,N),B(M,N)
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
CALL PARA_RANGE
& (1, M, NPROCS, MYRANK, ISTA, IEND)
ISTA2 = ISTA
IEND1 = IEND
IF (MYRANK == 0 ) ISTA2 = 2
IF (MYRANK == NPROCS-1) IEND1 = M-1
IUP = MYRANK + 1
IDOWN = MYRANK - 1
IF (MYRANK==NPROCS-1) IUP =MPI_PROC_NULL④
IF (MYRANK==0) IDOWN=MPI_PROC_NULL⑫
CALL PARA_TYPE_BLOCK2A
& (1, M, 1, N, MPI_REAL8, ITYPE) ⑫

```

図5-1-4(1)

```

DO J = 1, N
DO I = ISTA, IEND
A(I,J) = I + 10.0*J
ENDDO
ENDDO
CALL MPI_ISEND(A(IEND,1), 1, ITYPE,
& IUP, 1, MPI_COMM_WORLD, ISEND1, IERR) ⑥
CALL MPI_ISEND(A(ISTA,1), 1, ITYPE,
& IDOWN, 1, MPI_COMM_WORLD, ISEND2, IERR) ⑦
CALL MPI_IRECV(A(ISTA-1,1), 1, ITYPE,
& IDOWN, 1, MPI_COMM_WORLD, IRECV1, IERR) ⑧
CALL MPI_IRECV(A(IEND+1,1), 1, ITYPE,
& IUP, 1, MPI_COMM_WORLD, IRECV2, IERR) ⑨
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR) ⑩
DO J = 2, N-1
DO I = ISTA2, IEND1
B(I,J) = A(I-1,J) + A(I,J-1)
+ A(I,J+1) + A(I+1,J)
ENDDO
ENDDO
CALL MPI_FINALIZE(IERR)
END

```

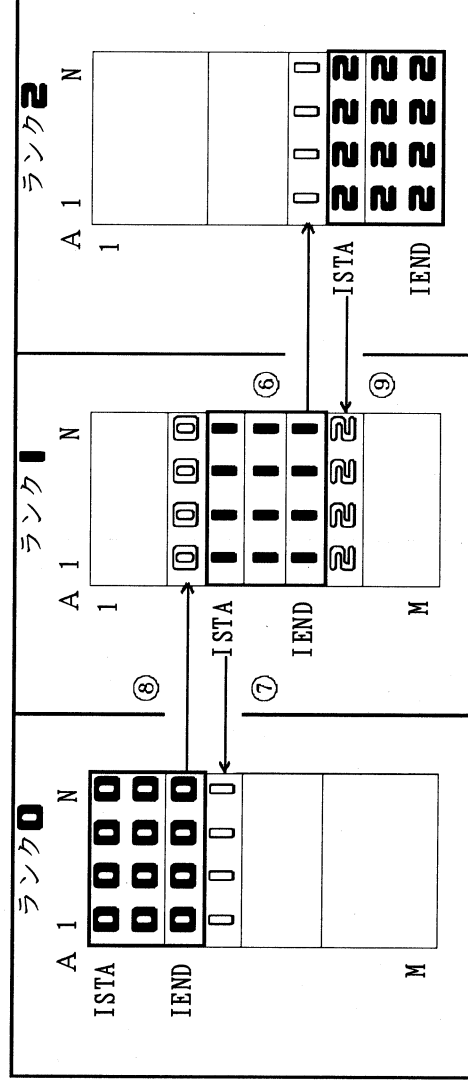


図5-1-4(2)

### 5-1-3 1, 2次元目でブロック分割した場合

配列A, Bを図5-1-2(3)のように1次元目と2次元目でブロック分割する場合(4-5-8節参照)のシフト方法を説明します(紙面の関係で配列A, Bの大きさは12×9とします)。

まず配列Aのデータの動きを図5-1-5(5)で概観します。例えばランク4のプロセスは、中央の図に示すように、上下左右の隣接する各プロセスとの間でシフトを行います。このうち左右のシフトは配列を2次元目でブロック分割した5-1-1節のシフトと同じ動作に、上下のシフトは配列を1次元目でブロック分割した5-1-2節のシフトと同じ動作になります。

並列化したプログラムを図5-1-5(4)に示し、以下で説明します。図5-1-5(4)(5)内の番号はそれぞれ対応しています。なお、図5-1-5(4)で⑥と⑦のように異なる記号の同じ数字がある場合、⑥が左右方向のシフト、⑦が上下方向のシフトに関連する部分を示します。

- 図5-1-5(4)の⑬で、配列A, Bを分割したときの1次元方向と2次元方向のプロセス数をIPROCS, JPROCSに設定します。この例では図5-1-5(5)に示すようにIPROCS, JPROCSともに3で、全部で9プロセスを使用します。⑭でIPROCS×JPROCSがNPROCSと等しいかどうか念のためチェックします。

- ⑮で図5-1-5(1)に示す配列ITABLE(以下プロセス格子テーブル)を作成します。これは各プロセスが自分の上下左右の隣のプロセスのランク値を知るために使用します。本例は非循環シフトなので、プロセスの存在しない×の部分にはMPI\_PROC\_NULL(4-6-2節参照)を入れます。また自プロセスのプロセス格子テーブル内の1次元目と2次元目の座標をMYRANKIとMYRANKJに設定します(例えば自分がランク1の場合、**◆**と**▼**に示すようにMYRANKI=0, MYRANKJ=1となります)。なお、以下ではテーブル内のプロセス値の順番を横方向に**0, 1, 2, …**としましたが、縦方向に**0, 1, 2, …**としてもかまいません。

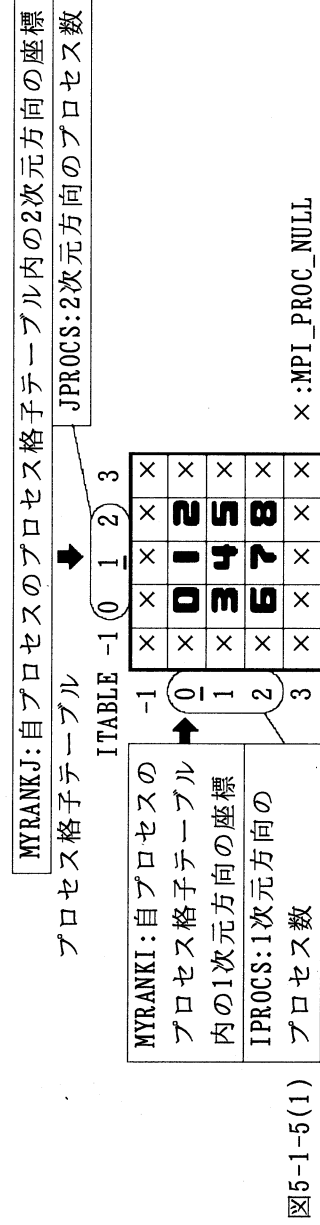


図5-1-5(1)

- ⑯と⑰で、自プロセスの配列A, Bの1次元目と2次元目の担当範囲を求めます。前節までの例では1つの次元のみで分割していたので、サブルーチンPARAM\_RANGEの引数を図5-1-5(2)の「今までの分割」のように指定しましたが、今回は2つの次元でそれぞれ分割するため、図5-1-5(2)の「1次元目の分割」と「2次元目の分割」のように指定します。

次に⑱と㉑で、図5-1-5(1)のプロセス格子テーブルを使用して、自プロセスの上下左右のプロセスのランク値を取得します。例えばランク4の場合、図5-1-5(3)のようになります。

- ⑳で、各プロセスが通信する境界の1列の要素数ILENを求めます(この値は⑯～㉑で使用します)。また㉒で、各プロセスが通信する境界の1行を表す派生データ型ITYPEを作成します(詳細は5-1-2節参照)。

	PARAM_RANGEで指定する引数	PARAM_RANGEで戻る引数
今までの分割	プロセス数	自分のランク
1次元目の分割	NPROCS	MYRANK
2次元目の分割	IPROCS	MYRANKI
	JPROCS	MYRANKJ

図5-1-5(2)

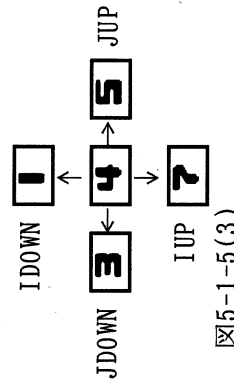


図5-1-5(3)

- 図5-1-1(1)①のループを図5-1-5(4)⑤でIとJの両方向で並列化します。次に配列Aの境界のデータのシフトを行います。図5-1-5(5)の左右方向のシフトは⑥～⑩で行い、上下方向のシフトは⑪～⑮で行います。このとき図5-1-5(3)の各変数を送信元または宛先プロセスのランク値として使用します。シフト終了後、図5-1-1(1)②のループを、図5-1-5(4)①でIとJの両方向で並列化します。

```

PROGRAM MAIN
IMPLICIT REAL*8(A-H,0-Z)
INCLUDE 'mpif.h'
PARAMETER (M=12,N=9)
DIMENSION A(M,N),B(M,N)
INTEGER ISTATUS(MPI_STATUS_SIZE)
PARAMETER (IPROCS=3, JPROCS=3)
INTEGER ITABLE(-1:IPROCS, -1:JPROCS)

CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
& (MPI_COMM_WORLD, MYRANK, IERR)

IF (NPROCS /= IPROCS*JPROCS) THEN
  IF (MYRANK==0) PRINT *, '====ERROR===='
  CALL MPI_FINALIZE(IERR)
  STOP
ENDIF
DO J = -1, JPROCS
  DO I = -1, IPROCS
    ITABLE(I, J) = MPI_PROC_NULL
  ENDDO
ENDDO
IRANK = 0
DO I = 0, IPROCS-1
  DO J = 0, JPROCS-1
    ITABLE(I, J) = IRANK
    IF (MYRANK == IRANK) THEN
      MYRANKI = I
      MYRANKJ = J
    ENDIF
    IRANK = IRANK + 1
  ENDDO
ENDDO
CALL PARA_RANGE
& (1, N, JPROCS, MYRANKJ, JSTA, JEND)
JSTA2 = JSTA
JEND1 = JEND
IF (MYRANKJ == 0 ) JSTA2 = 2
IF (MYRANKJ == JPROCS-1) JEND1 = N-1
CALL PARA_RANGE
& (1, M, IPROCS, MYRANKI, ISTA, IEND)
ISTA2 = ISTA
IEND1 = IEND
IF (MYRANKI == 0 ) ISTA2 = 2
IF (MYRANKI == IPROCS-1) IEND1 = M-1
JUP = ITABLE(MYRANKI , MYRANKJ+1)
JDOWN = ITABLE(MYRANKI , MYRANKJ-1)
IUP = ITABLE(MYRANKI+1, MYRANKJ )
IDOWN = ITABLE(MYRANKI-1, MYRANKJ )
ILEN = IEND-ISTA+1

```

```

CALL PARA_TYPE_BLOCK2A
& (1, M, 1, JEND-JSTA+1, MPI_REAL8, ITYPE)
DO J = JSTA, JEND
  DO I = ISTA, IEND
    A(I, J) = I + 10.0*J
  ENDDO
ENDDO
CALL MPI_ISEND(A(ISTA, JEND), ILEN,
& MPI_REAL8, JUP , 1,
& MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(A(ISTA, JSTA), ILEN,
& MPI_REAL8, JDOWN, 1,
& MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_ISEND(A(IEND, JSTA), 1,
& ITYPE, IUP , 1,
& MPI_COMM_WORLD, JSEND1, IERR)
CALL MPI_ISEND(A(ISTA, JSTA), 1,
& ITYPE, IDOWN, 1,
& MPI_COMM_WORLD, JSEND2, IERR)
CALL MPI_IRECV(A(ISTA, JSTA-1), ILEN,
& MPI_REAL8, JDOWN, 1,
& MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_IRECV(A(ISTA, JEND+1), ILEN,
& MPI_REAL8, JUP , 1,
& MPI_COMM_WORLD, IRECV2, IERR)
CALL MPI_IRECV(A(ISTA-1, JSTA), 1,
& ITYPE, IDOWN, 1,
& MPI_COMM_WORLD, JRECV1, IERR)
CALL MPI_IRECV(A(IEND+1, JSTA), 1,
& ITYPE, IUP , 1,
& MPI_COMM_WORLD, JRECV2, IERR)
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(JSEND1, ISTATUS, IERR)
CALL MPI_WAIT(JSEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
CALL MPI_WAIT(JRECV1, ISTATUS, IERR)
CALL MPI_WAIT(JRECV2, ISTATUS, IERR)
DO J = JSTA2, JEND1
  DO I = ISTA2, IEND1
    B(I, J) = A(I-1, J) + A(I, J-1)
    + A(I, J+1) + A(I+1, J)
  ENDDO
ENDDO
CALL MPI_FINALIZE(IERR)
END

```

ⓧ5-1-5(4)

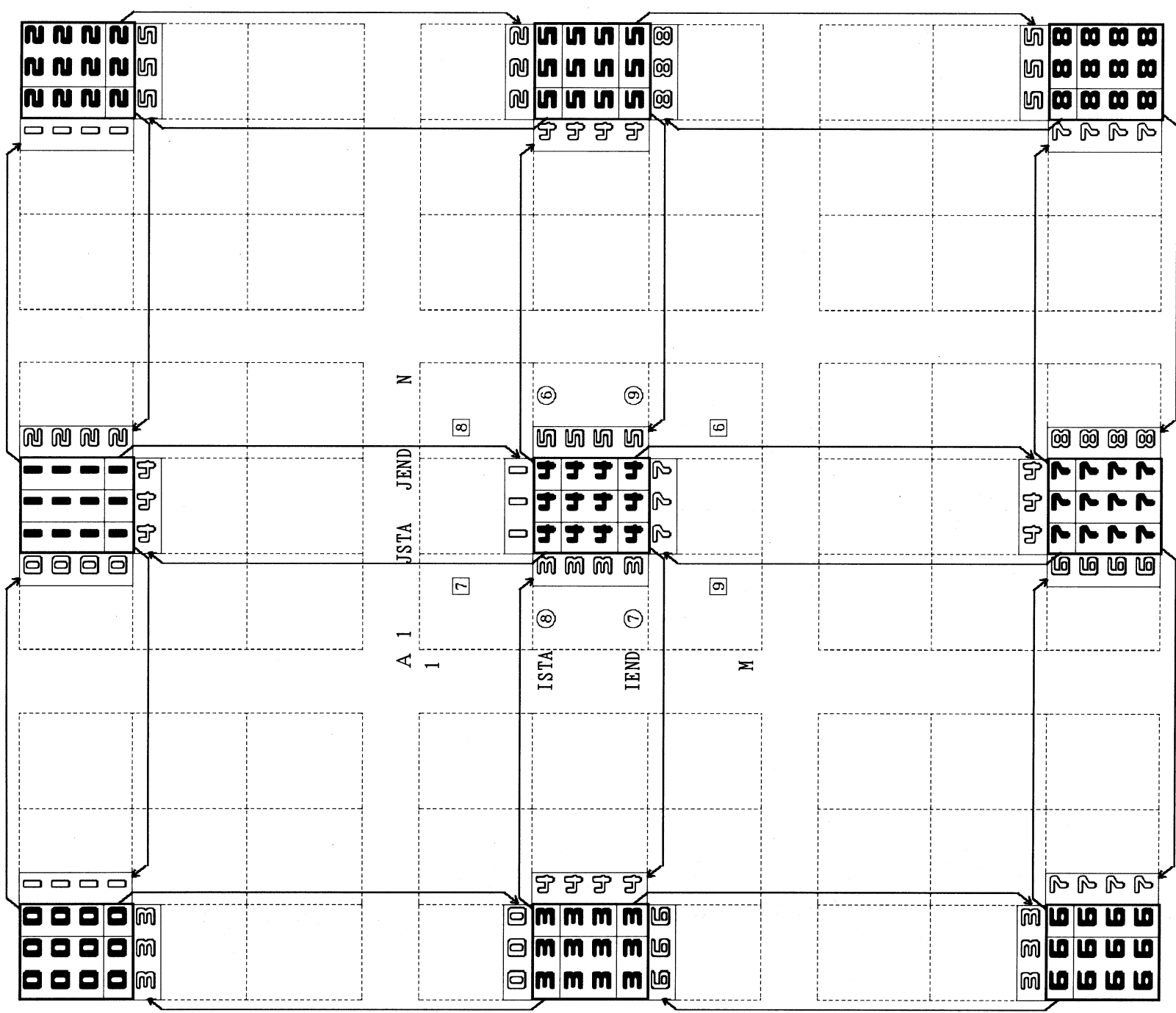


図5-1-5(5)

## 5-1-4 1, 2次元目でブロック分割した場合(斜めの要素を参照)

差分法のプログラムによっては、図5-1-6(1)の下線部が追加されることがあります。この場合、まず①で配列Aに適当な値を設定し(図5-1-6(2)の○, ●, ⊙)、②で配列Aの上下左右(●)と斜め隣の値(⊙)から配列Bの値(■)を求めます(紙面の関係で、図5-1-6(2)の配列の大きさは4×9になっています)。

```

PROGRAM MAIN
  IMPLICIT REAL*8(A-H, O-Z)
  PARAMETER (M=12, N=9)
  DIMENSION A(M, N), B(M, N)

  DO J = 1, N
    DO I = 1, M
      A(I, J) = I + 10.0 * J
    ENDDO
  ENDDO

  DO J = 2, N-1
    DO I = 2, M-1
      B(I, J) = A(I-1, J) + A(I, J-1)
             + A(I, J+1) + A(I+1, J)
             + A(I-1, J-1) + A(I+1, J-1)
             + A(I-1, J+1) + A(I+1, J+1)
    ENDDO
  ENDDO
END

```

図5-1-6(1)

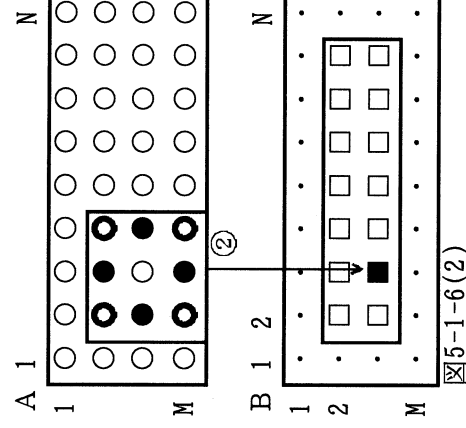


図5-1-6(2)

このプログラムを1次元目または2次元目のどちらから一方でブロック分割して並列化する場合には、5-1-1節と5-1-2節の方法をそのまま使用することができます。しかし1次元目と2次元目の両方でブロック分割する場合、5-1-3節の方法をそのまま使用することはできません。

その理由を説明します。並列化した場合の配列Aのデータの動きは図5-1-6(4)のようになりますが、中央の図に示すように、例えばリンク4のプロセスは図5-1-6(1)②の下線部の計算を行うため、①, ③, ⑤, ⑦の各要素の他に、四隅の②, ④, ⑥, ⑧の要素(斜め隣の各プロセスが持っている)が必要になります。ところが5-1-3節の方法では、図5-1-5(5)に示すようにこれらの要素を通信していませんので、この方法をそのまま使用することはできません。

四隅の要素は斜め隣のプロセスから直接送信することもできませんが、ここでは通信回数を減らすため直接通信しない方法を説明します。

まずデータの動きを図5-1-6(4)で概観します。最初に⑥~⑨で境界の1列を横方向に通信します。これを矢印(実線)で示します。この通信は図5-1-5(5)の通信と同じです。

横方向の通信が完了した後、⑥~⑨で境界の1行を縦方向に通信します。これを矢印(点線)で示します。このとき図5-1-5(5)の場合よりも左右に1要素分余分に通信します。例えばリンク1のプロセスは『②, ③, ④』をリンク4のプロセスに送信します(★の部分の通信)。

これによって、四隅の要素を含む全ての要素の通信が完了します。例えばリンク4のプロセスの四隅には、②, ④, ⑥, ⑧の要素が入っていますが、そのうち左上の要素②は、リンク②から直接リンク4に送られたのではなく、リンク②からリンク1を経由してリンク4に送られたことに注意して下さい。このような動作を行うためには、先に横方向の通信を行って、それが完了した後で縦方向の通信を(左右に1要素分余分に)行う必要があります。

次にプログラム例を図5-1-6(3)に示します。図5-1-6(3)(4)内の番号はそれぞれ対応しています。まず②で、縦方向に通信する1行(左右に余分な要素を1要素分持つ)の下限と上限をJJSTA, JJENDとし(図5-1-6(4)の↓参照)、この1行の派生データ型ITYPEを、3-5-3-2節で説明した自作のサブルーチンPARAM\_TYPE\_BLOCK2Aを使用して作成します(MPI-2が使用できるマシン環境の場合は、MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照)を使用できますが、派生データ型の起点が異なるので注意して下さい)。そして図5-1-6(1)⑤の計算が終了した後、まず⑥~⑨で境界の1列を横方向に通信し、それが完了した後、⑥~⑨で境界の1行を派生データ型ITYPEを使用して縦方向に通信します。



```

PROGRAM MAIN
IMPLICIT REAL*8(A-H,0-Z)
INCLUDE 'mpif.h'
PARAMETER (M=12,N=9)
DIMENSION A(M,N),B(M,N)
INTEGER ISTATUS(MPI_STATUS_SIZE)
PARAMETER(IPROCS=3,JPROCS=3)
INTEGER ITABLE(-1:IPROCS,-1:JPROCS)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD,NPROCS,IERR)
& (MPI_COMM_WORLD,MYRANK,IERR)
& (MPI_COMM_WORLD,MYRANK,IERR)
IF (NPROCS /= IPROCS*JPROCS) THEN
IF (MYRANK==0) PRINT *, '===ERROR==='
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
DO J = -1, JPROCS
DO I = -1, IPROCS
ITABLE(I,J) = MPI_PROC_NULL
ENDDO
ENDDO
IRANK = 0
DO I = 0, IPROCS-1
DO J = 0, JPROCS-1
ITABLE(I,J) = IRANK
IF (MYRANK == IRANK) THEN
MYRANKI = I
MYRANKJ = J
ENDIF
IRANK = IRANK + 1
ENDDO
ENDDO
CALL PARA_RANGE
& (1,N,JPROCS,MYRANKJ,JSTA,IEND)
JSTA2 = JSTA
JEND1 = JEND
IF (MYRANKJ == 0 ) JSTA2 = 2
IF (MYRANKJ == JPROCS-1) JEND1 = N-1
CALL PARA_RANGE
& (1,M,IPROCS,MYRANKI,ISTA,IEND)
ISTA2 = ISTA
IEND1 = IEND
IF (MYRANKI == 0 ) ISTA2 = 2
IF (MYRANKI == IPROCS-1) IEND1 = M-1
JUP = ITABLE(MYRANKI ,MYRANKJ+1)
JDOWN = ITABLE(MYRANKI ,MYRANKJ-1)
IUP = ITABLE(MYRANKI+1,MYRANKJ )
IDOWN = ITABLE(MYRANKI-1,MYRANKJ )
ILEN = IEND-ISTA+1

```

```

JSTA = MAX(1,JSTA-1)
JJEND = MIN(N,JEND+1)
CALL PARA_TYPE_BLOCK2A
& (1,M,1,JJEND-JJSTA+1,MPI_REAL8,ITYPE)
DO J = JSTA, JEND
DO I = ISTA, IEND
A(I,J) = I + 10.0*J
ENDDO
ENDDO
CALL MPI_ISEND(A(ISTA,JEND),ILEN,
& MPI_REAL8,JUP ,1,
& MPI_COMM_WORLD,ISEND1,IERR)
CALL MPI_ISEND(A(ISTA,JSTA),ILEN,
& MPI_REAL8,JDOWN,1,
& MPI_COMM_WORLD,ISEND2,IERR)
CALL MPI_IRECV(A(ISTA,JSTA-1),ILEN,
& MPI_REAL8,JDOWN,1,
& MPI_COMM_WORLD,IRecv1,IERR)
CALL MPI_IRECV(A(ISTA,JEND+1),ILEN,
& MPI_REAL8,JUP ,1,
& MPI_COMM_WORLD,IRecv2,IERR)
CALL MPI_WAIT(ISEND1,ISTATUS,IERR)
CALL MPI_WAIT(ISEND2,ISTATUS,IERR)
CALL MPI_WAIT(IRecv1,ISTATUS,IERR)
CALL MPI_WAIT(IRecv2,ISTATUS,IERR)
CALL MPI_ISEND(A(IEND,JJSTA),1,
& ITYPE,IUP ,1,
& MPI_COMM_WORLD,JSEND1,IERR)
CALL MPI_ISEND(A(ISTA,JJSTA),1,
& ITYPE,IDOWN,1,
& MPI_COMM_WORLD,JSEND2,IERR)
CALL MPI_IRECV(A(ISTA-1,JJSTA),1,
& ITYPE,IDOWN,1,
& MPI_COMM_WORLD,JRECV1,IERR)
CALL MPI_IRECV(A(IEND+1,JJSTA),1,
& ITYPE,IUP ,1,
& MPI_COMM_WORLD,JRECV2,IERR)
CALL MPI_WAIT(JSEND1,ISTATUS,IERR)
CALL MPI_WAIT(JSEND2,ISTATUS,IERR)
CALL MPI_WAIT(JRECV1,ISTATUS,IERR)
CALL MPI_WAIT(JRECV2,ISTATUS,IERR)
DO J = JSTA2, JEND1
DO I = ISTA2, IEND1
B(I,J) = A(I-1,J ) + A(I ,J-1)
+ A(I ,J+1) + A(I+1,J )
+ A(I-1,J-1) + A(I+1,J-1)
+ A(I-1,J+1) + A(I+1,J+1)
ENDDO
ENDDO
CALL MPI_FINALIZE(IERR)
END

```

図5-1-6(3)

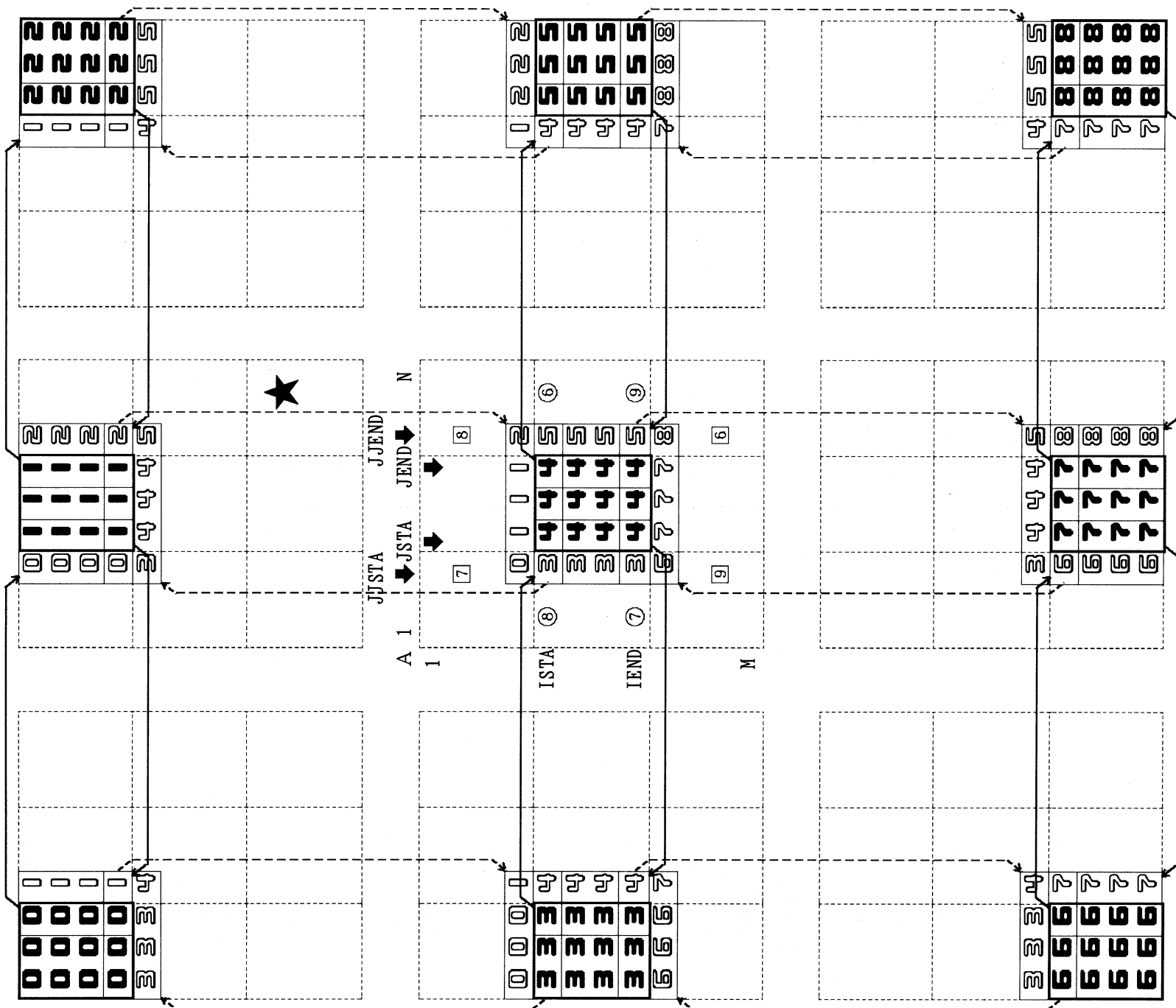


图 5-1-6(4)

### 5-1-5 配列を縮小する場合

本節では、4-5-7節で説明した動的割振りをを用いて配列を縮小する方法について説明します。前節までのいずれの分割でも縮小方法は同じなので、配列を2次元方向で分割した図5-1-3(1)(2)を例とし、縮小後のプログラムを図5-1-7(1)に、縮小した配列A,Bを図5-1-7(2)の実線部に示します。

図5-1-7(1)では、まず⑦で縮小する配列A,Bの動的割振りを宣言します。そして③で2次元方向の自分の担当範囲を求めた後、④で配列Aを確保します。図5-1-7(2)の上段のランク1に示すように、JSTAの前とJENDの後にシフト用の1列が必要になるので、J方向の大きさは(JSTA-1)~(JEND+1)となります。ただし、この例では端のプロセスであるランク0ではJSTA-1が、ランク2ではJEND+1が不要なので、組込関数MIN,MAXを使用し調整します。

次に⑤で配列Bを確保します。⑥のD0ループの反復から分かるように、1次元方向の大きさは2~M-1、2次元方向の大きさはJSTA2~JEND1となります。

なお、派生データ型(3-5節参照)を使用して通信した配列を縮小する場合、4-5-7-2節の「派生データ型を使用した通信の修正」で説明した修正が必要になります。

```

:
INCLUDE 'mpif.h'
PARAMETER (M=4,N=9)
REAL*8,ALLOCATABLE::A(:,,:),B(:,:)
CALL PARA_RANGE
& (1,N,NPROCS,MYRANK,JSTA,JEND)
JSTA2 = JSTA
JEND1 = JEND
IF (MYRANK == 0 ) JSTA2 = 2
IF (MYRANK == NPROCS-1) JEND1 = N-1
ALLOCATE(A(M,MAX(JSTA-1,1):MIN(JEND+1,N)))
ALLOCATE(B(2:M-1,JSTA2:JEND1))
:

```

図5-1-7(1)

```

D0 J = JSTA, JEND
D0 I = 1, M
A(I,J) = I + 10.0*J
ENDDO
ENDDO
:
D0 J = JSTA2, JEND1
D0 I = 2, M-1
B(I,J) = A(I-1,J) + A(I,J-1)
+ A(I,J+1) + A(I+1,J)
ENDDO
ENDDO
CALL MPI_FINALIZE(IERR)
:

```

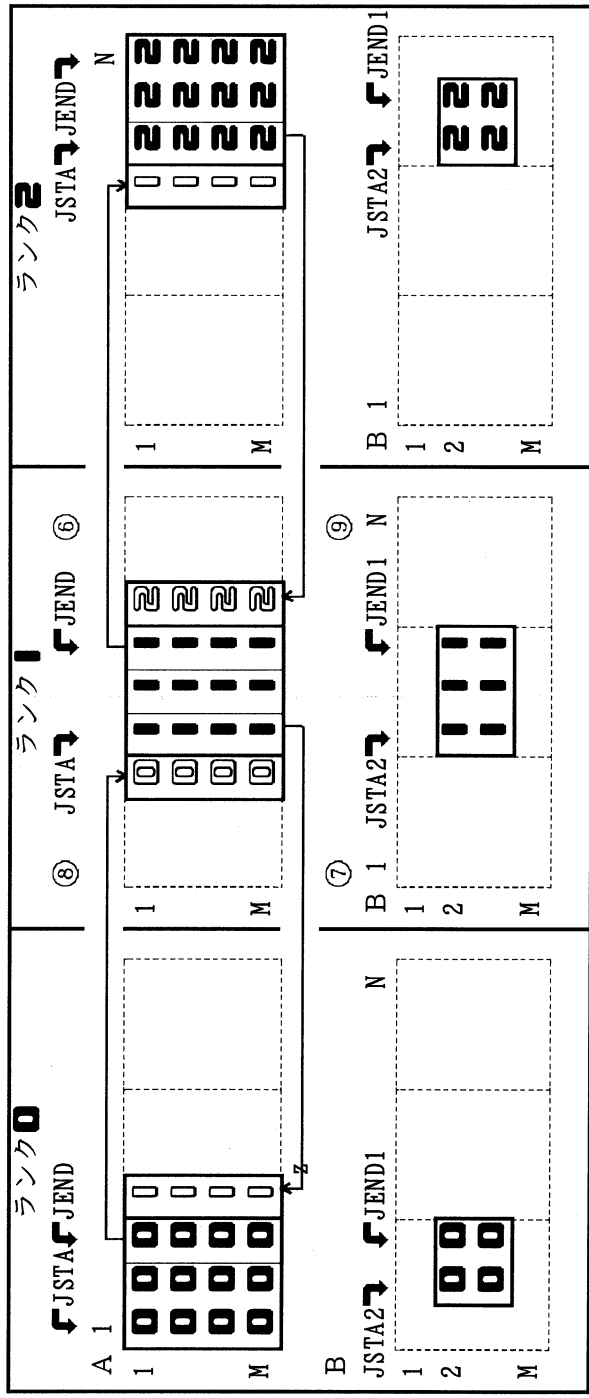


図5-1-7(2)

## 5-1-6 その他の考慮点

例えば図5-1-8(1)の①、②のD0ループをJで並列化する場合、配列B,D,Fに対してシフトを行なう必要が  
あります。差分法のプログラムを並列化する場合、このように通常多くの配列に対し、多くの箇所でシフトを  
行う必要がありますが、通信が必要になるたびに通信ルーチンを毎回記述するとプログラムがゴチャゴチャ  
になります。このような場合、図5-1-8(2)のような方法で並列化するとプログラムがすっきりします。

- 並列化に関係のある汎用的な変数は全て③のMODULE文内に記述します(詳細は4-8-1節を参照して下  
さい)。MODULE文は少し使い易い方にクセがあるので、4-5-7節を参照して下さい。
- 実際にシフトを行う部分を⑥のようにサブルーチン化します。
- ①と②のD0ループの直前の④と⑤で、シフトを行うサブルーチンSHIFTをコールします。

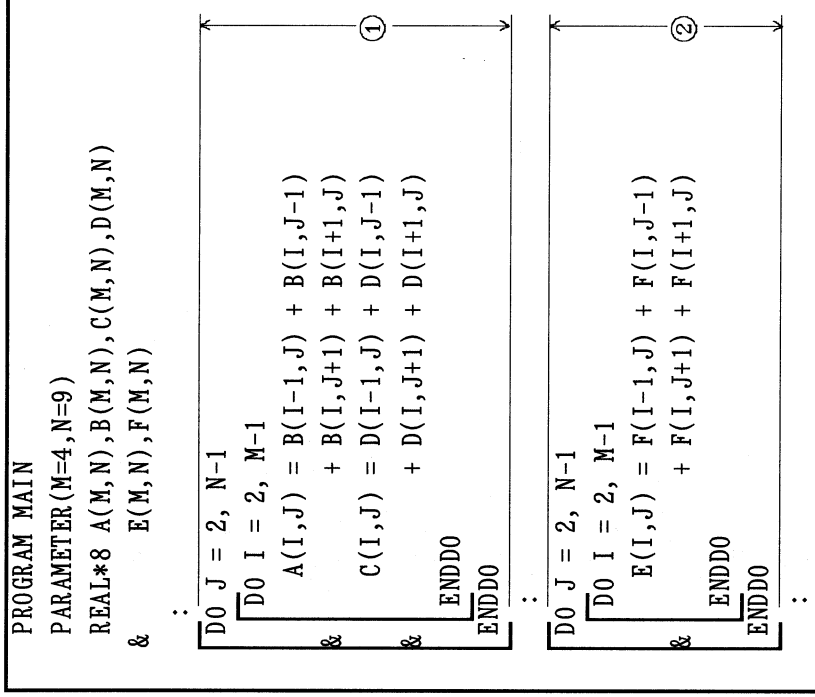


図5-1-8(1)

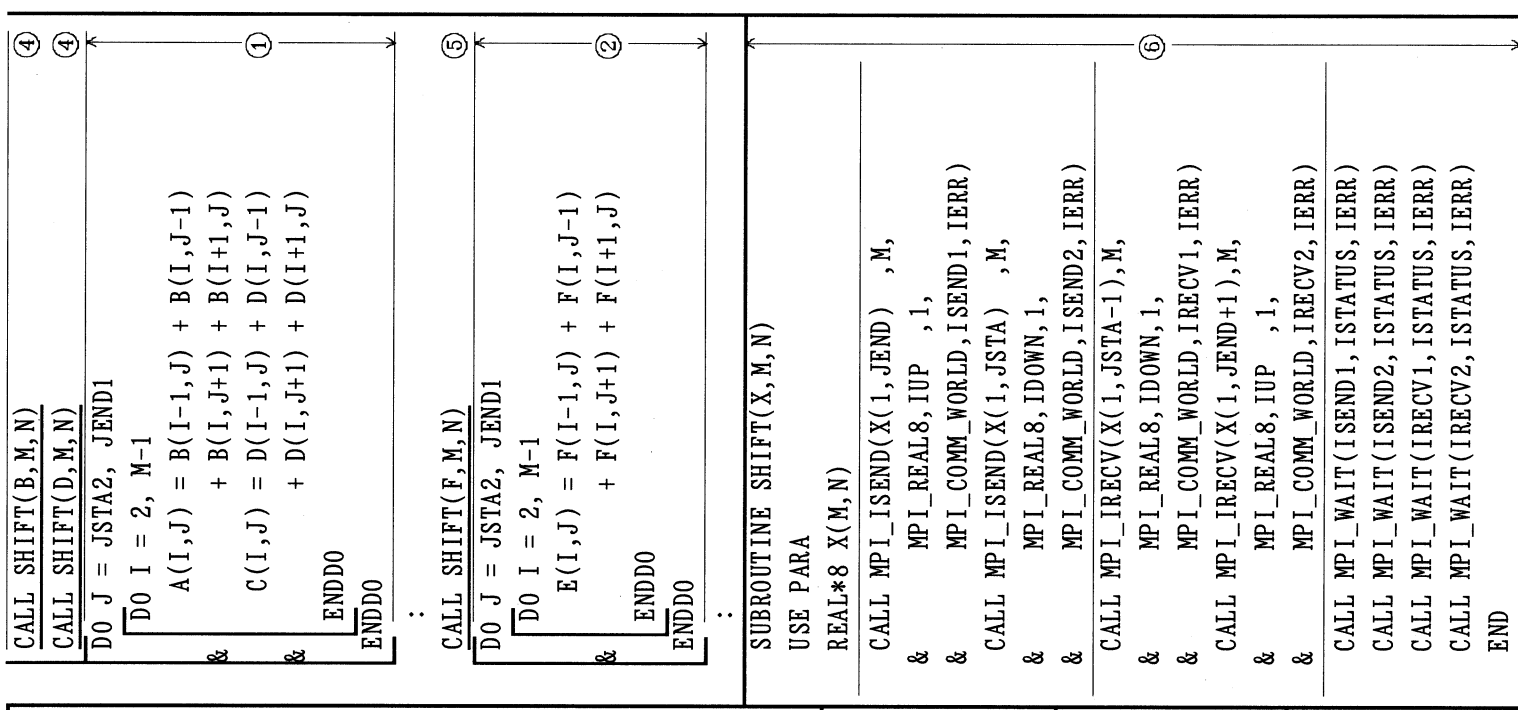


図5-1-8(2)

## 5-2-1 有限要素法の並列化

有限要素法は、要素の形状が不規則なため通信がしにくく、差分法と比べて一般に並列化が難しくなりません。まず単体版の有限要素法のプログラムをそのまま並列化する場合を想定します。

陰解法のプログラムの場合、不規則対称疎行列の連立一次方程式を解く部分がホットスポットになります。これを直接法で解くための、数値計算サブルーチンの並列版(第6章参照)が提供されているかどうかは未調査です。

一方陽解法のプログラムの場合、本節に述べるような方法で(一応)並列化することができます。

なお、市販の有限要素法プログラムの並列版では、領域分割法という方法で並列化するのが主流になっているようです。領域分割法では、例えば3ノードで計算する場合、計算領域を3つの完全に独立した領域として解き、領域の継ぎ目の部分を調整するという方法をとります。従って単体版のプログラムのそのまま並列化できるわけではなく、計算のアルゴリズム自体を並列用に変更する必要があります。領域分割法については参考文献[16],[17]などを参照して下さい。

### ■ オリジナル版プログラム

以下では有限要素法の陽解法モドキのプログラムである図5-2-2で説明します。配列Yは要素、配列Sは節点の配列です。計算領域は図5-2-1(1)のようになっています。□が要素番号、○が節点番号を示します。なお、本例では説明を簡単にするために四角形要素を使用し、計算領域の形状も単純ですが、三角形要素や領域の形状が複雑な場合でも同じ方法が適用できます。

以下に図5-2-2のプログラムの動作を示します。

- [1]で配列Y, Sに初期値を設定します。
- [2]では図5-2-1(1)の矢印に示すように、要素の値を使用してその要素の四隅の節点の値を更新します。このとき、要素番号と周囲の4節点の節点番号の対応表である配列INDEX(図5-2-1(5))を使用します。
- [3]では図5-2-1(2)の▼に示すように、節点の値を更新します。
- [4]では図5-2-1(3)の矢印に示すように、要素の値を、その要素の四隅の節点の値を使用して更新します。
- [5]では図5-2-1(4)の◆に示すように、要素の値を更新します。

一般に有限要素法(陽解法)のプログラムでは、このように3つのタイプのループ(要素のループ、節点のループ、要素↔節点のループ)が現れるのが特徴です。

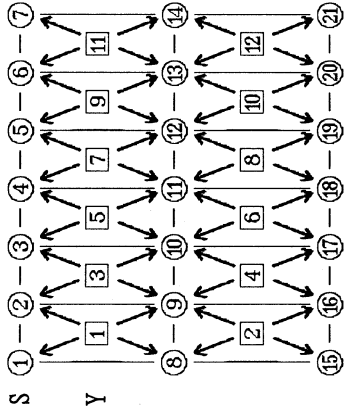


図5-2-1(1)

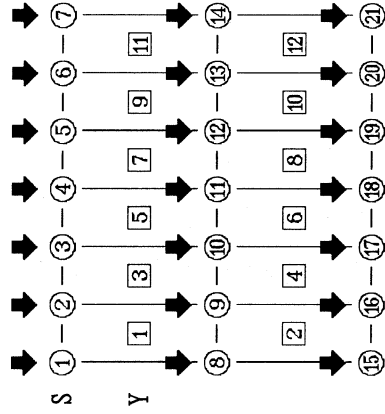


図5-2-1(2)

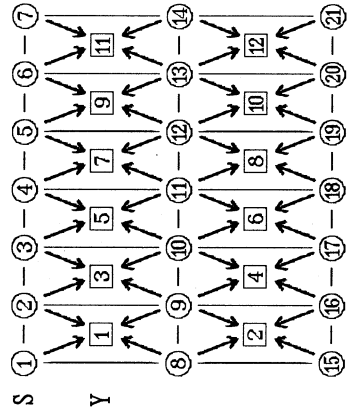


図5-2-1(3)

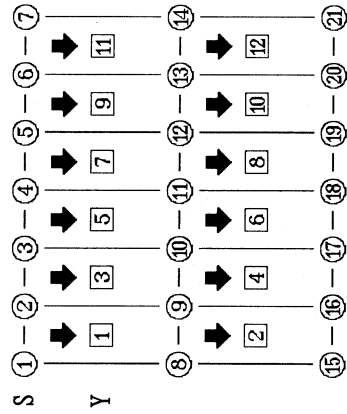


図5-2-1(4)

INDEX

1	2	3	4	5	6	7	8	9	10	11	12
①	②	③	④	⑤	⑥	⑦	⑧	⑨	⑩	⑪	⑫
⑬	⑭	⑮	⑯	⑰	⑱	⑲	⑳	㉑	㉒	㉓	㉔

図5-2-1(5)

```

:
PARAMETER(IYMAX=12, ISMAX=21)
REAL*8 Y(IYMAX), S(ISMAX)
INTEGER INDEX(4, IYMAX)
:
D0 IY=1, IYMAX
  Y(IY) = IY*10.0
ENDDO
D0 IS=1, ISMAX
  S(IS) = IS*100.0
ENDDO
D0 ITIME=1, 10
  D0 IY=1, IYMAX
    要素から4節点を更新
    D0 J=1, 4
      S(INDEX(J, IY))=S(INDEX(J, IY))+Y(IY)[2]
    ENDDO
  ENDDO
  D0 IS=1, ISMAX
    節点を更新
    S(IS) = S(IS)*0.25
  ENDDO
  D0 IY=1, IYMAX
    4節点から要素を更新
    D0 J=1, 4
      Y(IY) = Y(IY) + S(INDEX(J, IY))
    ENDDO
  ENDDO
  D0 IY=1, IYMAX
    要素を更新
    Y(IY) = Y(IY)*0.25
  ENDDO
ENDDO
PRINT *, '結果', S, Y
:

```

図5-2-2

## 5-2-2 並列化の方法1

まず、4-6-6節で説明した『重ね合せ』を用いて簡単に(手抜きで)並列化する方法について説明します。並列化したプログラムを図5-2-5に示します。

●[1]で、要素数IYMAXをブロック分割し、各プロセスの担当する要素の範囲IYSTA, IYENDを求めます(図5-2-3(1)参照)。同様に節点数ISMAYをブロック分割し、各プロセスの担当する節点の範囲ISSTA, ISENDを求めます(図5-2-3(2)参照)。要素の分割と節点の分割には何の関係もないことに注意して下さい。

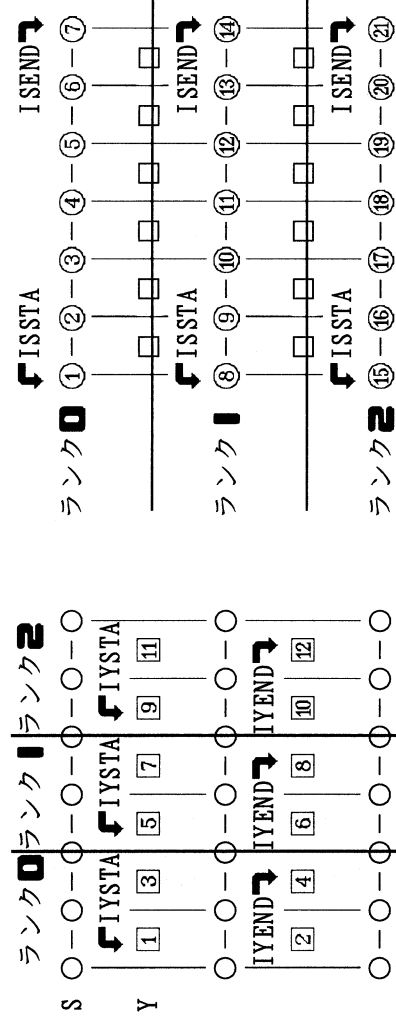


図5-2-3(1) 要素の分割

図5-2-3(2) 節点の分割

- [2]で配列Y,Sをゼロクリアしておきます(この目的については後述します)。次に[3]を並列化します。
- [5](図5-2-2の[2]の部分)を並列化します。配列Sには現在までの計算結果の累計が入っています。ここでは配列Sを直接更新するのではなく、一時配列SSに現タイムステップでの増分を求め、それを配列Sに加算するという方法をとります。  
まず[4]で一時配列SSを初期化し、次に[5]でこのタイムステップでの増分を求めます。[5]が終了すると図5-2-4(1)のようになります。プロセス間の境界にある節点には、完全な増分値ではなく各プロセスが計算した部分増分が求められます。例えば③の場合、ランク0のプロセスには部分増分③が、ランク1のプロセスには部分増分③と③の値は合計して完全な増分③にしてから配列Sに加算する必要があります。しかしどの節点に部分増分が入っているのか(つまりどの節点がプロセス間の境界にあるのか)を調べるのが面倒なので、ここでは手抜きの方法として、[6]で配列SSのすべての要素を重ね合わせて結果を全プロセスの一時配列SSSに求め(図5-2-4(2))、それを[7]で元の配列Sに加算するようにしました。なお[4]で配列SSをゼロクリアしたのは、[6]で配列SSについて重ね合せを行うためです。

- [8](図5-2-2の[3]の部分)を並列化します。この結果、図5-2-4(3)で↓の付いた節点が更新されます。
- [10](図5-2-2の[4]の部分)を並列化します。図5-2-4(5)に示すように、例えばランク0のプロセスは①~③、⑨~⑫、⑮~⑰を使用して計算しますが、[6]が終了した直後には図5-2-4(3)の↓に示すように①~⑷しか持っていないません。そこでまず[9]で配列Sの値を重ね合わせて結果を全プロセスの一時配列SSに入れ(図5-2-4(4))、その後[10]でSSを参照するようにします(図5-2-4(5))。なお[2]で配列Sをゼロクリアしたのは、[9]で配列Sについて重ね合せを行うためです。
- [11](図5-2-2の[5]の部分)を並列化します。この結果、図5-2-4(6)の↓の付いた要素が更新されます。
- タイムステップループが終了した後、[2]で、各プロセスが計算した配列YとSの値を重ね合せによってランク0のプロセスに集めます。配列Y,Sについて重ね合せを行うため、[2]で配列Y,Sをゼロクリアしてあります。

この方法は、次に述べる必要最小限のデータだけを通信する方法と比べて修正が簡単ですが、重ね合せによって余分な通信を行なうため通信時間が増大します。ただし計算部分が多いプログラムの場合はこの方法でも十分『ペイ』できます。また、最終的には次に述べる方法を使用する場合でも、とりあえず並列化がうまくいくかを確認したい場合に上記の方法を使用するという手もあります。

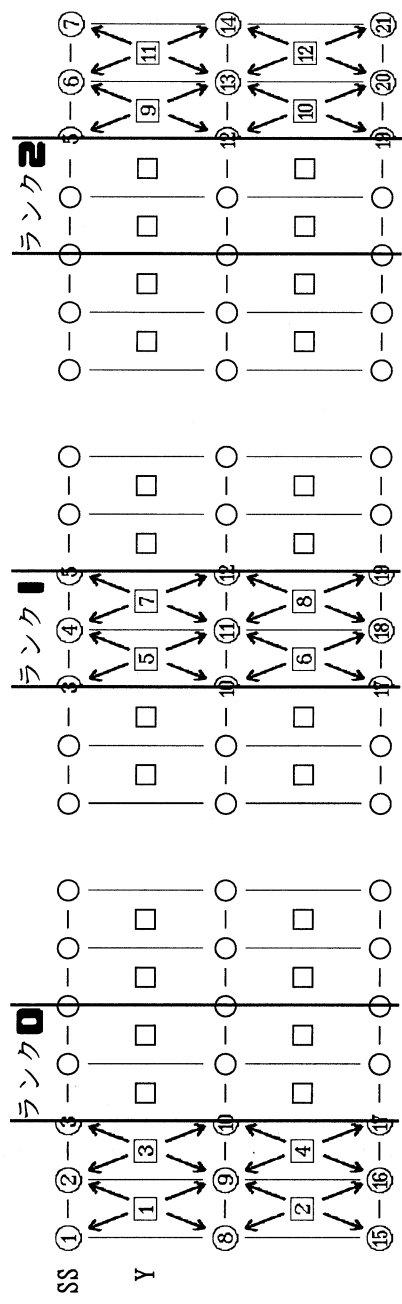


図5-2-4(1)

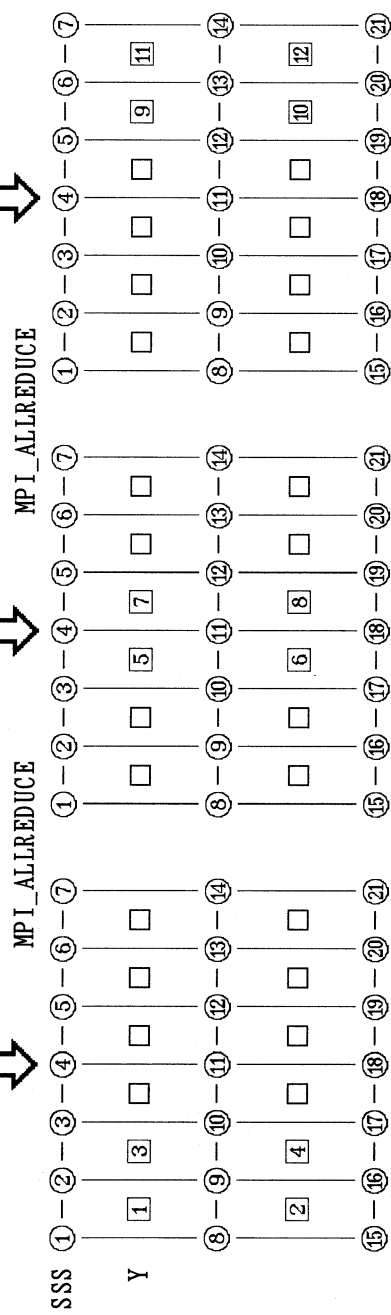


図5-2-4(2)

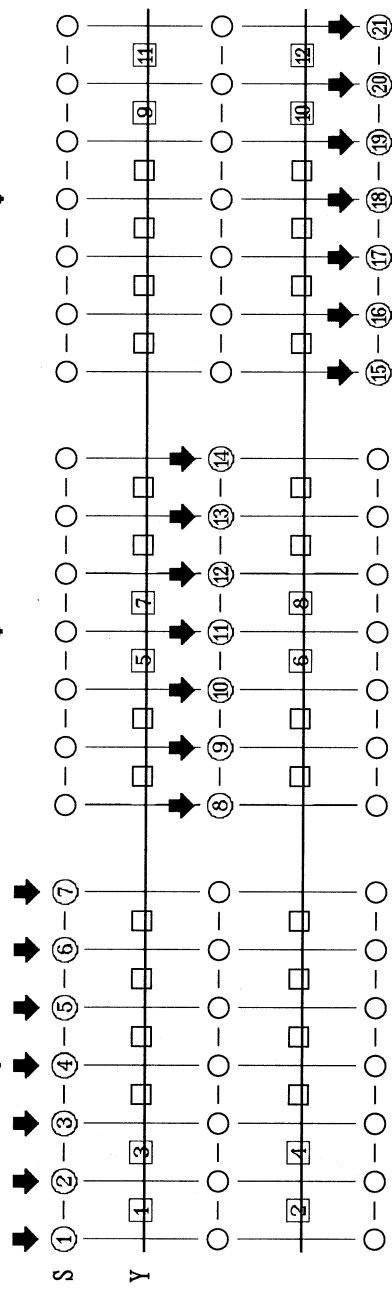


図5-2-4(3)

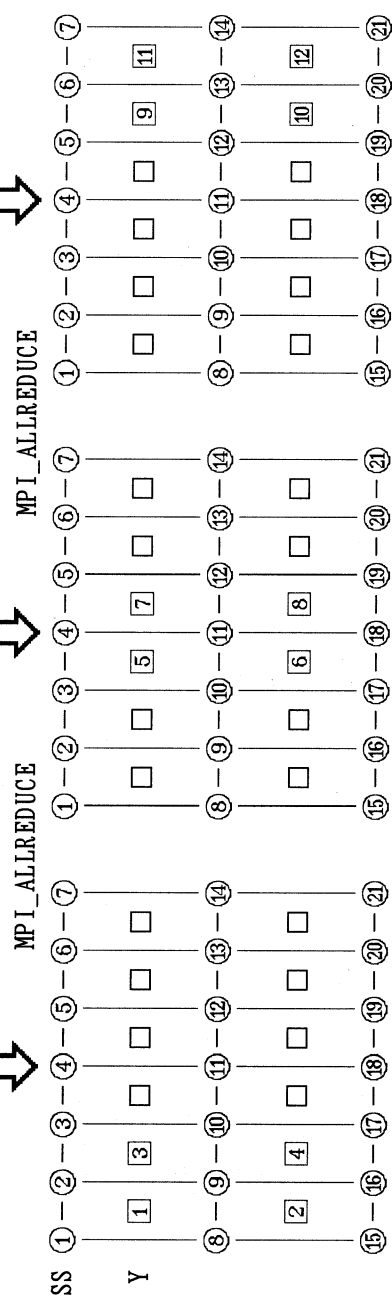


図5-2-4(4)



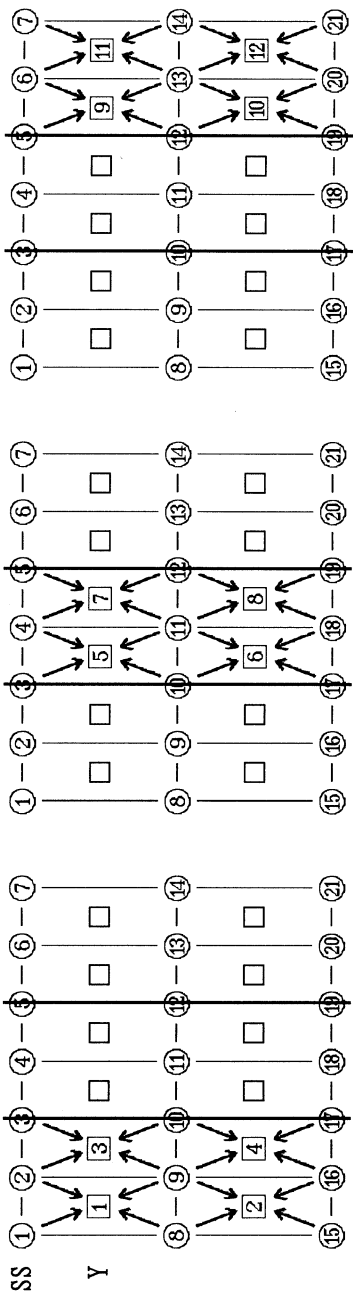


図5-2-4(5)

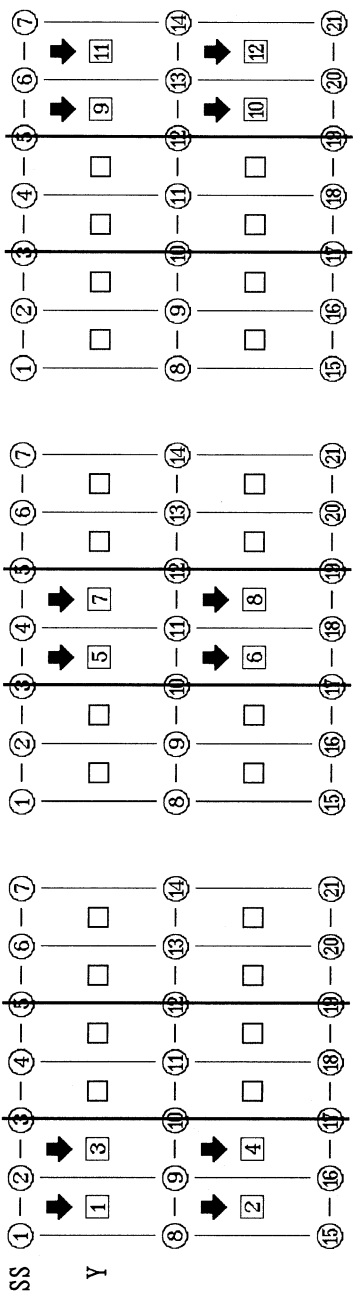


図5-2-4(6)

```

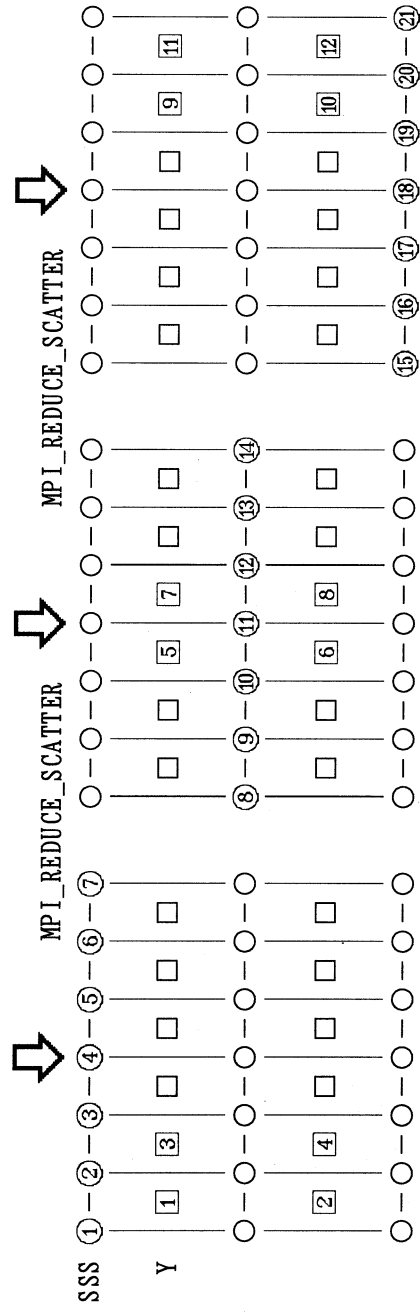
:
INCLUDE 'mpif.h'
PARAMETER (IYMAX=12, ISMAX=21)
REAL*8 Y(IYMAX), S(ISMAX)
REAL*8 YY(IYMAX), SS(ISMAX), SSS(ISMAX)
INTEGER INDEX(4, IYMAX)
:
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(~, NPROCS, IERR)
CALL MPI_COMM_RANK(~, MYRANK, IERR)
CALL PARA_RANGE
& (1, IYMAX, NPROCS, MYRANK, IYSTA, IYEND) [1]
CALL PARA_RANGE
& (1, ISMAX, NPROCS, MYRANK, ISSTA, ISEND)
DO IY=1, IYMAX
  Y(IY) = 0.0
ENDDO
DO IS=1, ISMAX
  S(IS) = 0.0
ENDDO
DO IY=IYSTA, IYEND
  Y(IY) = IY*10.00
ENDDO
DO IS=ISSTA, ISEND
  S(IS) = IS*100.0
ENDDO
DO ITIME=1, 10
  DO IS=1, ISMAX
    SS(IS) = 0.0
  ENDDO
  
```

図5-2-5

```

DO IY=IYSTA, IYEND 要素から4節点を更新
  DO J=1, 4
    SS(INDEX(J, IY))=SS(INDEX(J, IY))+Y(IY)[5]
  ENDDO
ENDDO
CALL MPI_ALLREDUCE(SS, SSS, ISMAX, MPI_REAL8,
  MPI_SUM, MPI_COMM_WORLD, IERR) [6]
DO IS=ISSTA, ISEND
  S(IS) = S(IS) + SSS(IS)
ENDDO
DO IS=ISSTA, ISEND 節点を更新
  S(IS) = S(IS)*0.25
ENDDO
CALL MPI_ALLREDUCE(S, SS, ISMAX, MPI_REAL8,
  MPI_SUM, MPI_COMM_WORLD, IERR) [9]
DO IY=IYSTA, IYEND 4節点から要素を更新
  DO J=1, 4
    Y(IY) = Y(IY) + SS(INDEX(J, IY)) [10]
  ENDDO
ENDDO
DO IY=IYSTA, IYEND 要素を更新
  Y(IY) = Y(IY)*0.25
ENDDO
ENDDO
CALL MPI_REDUCE(S, SS, ISMAX, MPI_REAL8,
  MPI_SUM, 0, MPI_COMM_WORLD, IERR)
&
CALL MPI_REDUCE(Y, YY, IYMAX, MPI_REAL8,
  MPI_SUM, 0, MPI_COMM_WORLD, IERR)
&
IF (MYRANK==0) PRINT *, '結果', S, Y
CALL MPI_FINALIZE(IERR)
:
  
```

余談ですが、上記の方法の通信部分には改良の余地があります。まず[6]で各プロセスに全ての節点の結果を渡しましたが、MPI\_REDUCE\_SCATTERを使用して、各プロセスごとに必要な節点の結果のみを送信すること  
もできます(4-6-6節参照)。このとき図5-2-4(2)は以下ようになります。



この方法を行うには図5-2-5を以下の③,⑥,⑨のように修正します。

また、[9]では『重ね合せ』で全ての節点の値を通信しましたが、図5-2-4(3)が終了した後で、各プロセスは、自分が計算した節点部分(↓の部分)のみを他の全プロセスに送信するようにすれば通信量が減ります。  
この方法を行うには図5-2-5を以下の③,⑥,⑨のように修正します。

```

:
PARAMETER (NCPU=3)
INTEGER ICNT(0:NCPU-1), IDISP(0:NCPU-1)
:
DO IRANK = 0, NPROCS-1
  CALL PARA_RANGE
  &      (1, ISMAX, NPROCS, IRANK, ISSTA, ISEND)
  ICNT(IRANK) = ISEND-ISSTA+1
  IDISP(IRANK) = ISSTA-1
ENDDO
:
CALL MPI_REDUCE_SCATTER(SS, SSS(ISSTA),
  ICNT, MPI_REAL8, MPI_SUM,
  &      MPI_COMM_WORLD, IERR)
:
CALL MPI_ALLGATHERV(S(ISSTA),
  &      ISEND-ISSTA+1, MPI_REAL8, SS, ICNT,
  &      IDISP, MPI_REAL8, MPI_COMM_WORLD, IERR)
:

```

↑ ③ : 宣言文に追加する。

↑ ⑥ : [1]の前に入れる。

↑ ⑨ : [6]を置き換える。

↑ ⑨ : [9]を置き換える。

前節で説明した方法は修正は比較的簡単ですが、全てのデータを「重ね合せ」で通信するため通信量が多くなります。本節では必要最小限のデータのみを通信して並列化する方法について説明します。

まず計算領域の分割方法について説明します。要素は図5-2-6(1)に示すように要素数でブロック分割します。次に節点ですが、前節の方法のように節点数でブロック分割するのではなく、各プロセスは自分が担当する要素の四隅の節点を担当します。ただし複数のプロセスの境界にある節点(図5-2-6(1)の○で囲んだ部分)は、ランクの最も小さいプロセスが担当することになります。例えば③,⑩,⑭は(ランク■でなく)ランク■が、⑤,⑯,⑲は(ランク■でなく)ランク■が担当します。まとめると、各プロセスが担当する節点番号は図5-2-6(2)のようになります。並列化したプログラムを図5-2-10に示し、以下で説明します。

### ■ テーブル類の作成

まず、後の並列計算で使用するテーブル類(図5-2-6(4)～図5-2-6(9))を作成します。

- [1]でテーブル類を初期化し、[2]で図5-2-6(4)に示すテーブルIDISPに値を設定します。IDISPは[4],[18],[23]の通信の引数を指定します(付録のMPI\_ALLTOALLVの引数sd displsとrd displs、MPI\_GATHERVの引数displs参照)。
- [3]で要素数IYMAXをブロック分割し(4-5-4節参照)、各プロセスが担当する要素の開始番号をIIYSTAに、担当する要素数をIYCNTにセットします(図5-2-6(5)参照)。
- [4]で作業配列ITEMPに値を設定します(図5-2-6(6)参照)。図5-2-6(1)で、例えばはランク■が担当する要素は①～④ですが、これらの要素の四隅にある節点番号を図5-2-6(3)の配列INDEXを使用して調べると、①～③,⑥～⑩,⑬～⑰でありますが、ここで配列ITEMP(～,■)のこれらの節点番号のところに「1」をセットします(ITEMPのその他の部分はゼロになっている必要があるため、[1]でゼロクリアしています)。
- [5]で改めて、自分が担当する要素の開始番号をIYSTAに、終了番号をIYENDに設定します(IYSTAは[3]で求めたIIYSTA(MYRANK)と同じ値です)。
- [6]でテーブルIBUF1, ICNT1, IBUF2, ICNT2を作成します。例えば図5-2-6(1)で節点③はランク■と■の境界にあります。この場合、前述のようにランク■と■のうちランクの小さいランク■のプロセスが節点③を担当します。これを[6]では次のように処理します。図5-2-6(6)の節点③(○で囲んだ部分)に複数の「1」が付いています。最も左(本例ではランク■)の「1」のみを残し、それ以外(本例ではランク■)の「1」を「0」に変更します(1→0の部分)。「1」を「0」に変更した場合は、ランク■のプロセスではIBUF2(～,■)に③を入れてICNT2(■)に1を加え(図5-2-6(9)の○で囲んだ部分)、ランク■のプロセスではIBUF1(～,■)に③を入れてICNT1(■)に1を加えます(図5-2-6(8)の○で囲んだ部分)。

図5-2-6(8)(9)から分かるように、作成されたIBUF1, ICNT1, IBUF2, ICNT2の値は、各プロセスで異なることに注意して下さい(それ以外のテーブル類の値は全プロセスで同一です)。

- [7]では先程作成したITEMPを使用して、各プロセスが担当する節点番号をISNUMに、節点の個数をISCNTにセットします(図5-2-6(7)参照)。以上で並列計算が必要となるテーブル類の作成は完了しました。

### ■ 並列計算部分

- [8]と[9]で図5-2-2の[1]を並列化します。このうち[9]は節点番号で並列化するので、先程作成したテーブルISNUMとISCNT(図5-2-6(7))を使用します。

- [10]からタイムステップのループが開始します。以後の説明では図5-2-7(1)～(4)も参照して下さい(図中で[1]などの数字は図5-2-10内の数字に対応します)。

図5-2-7(1)で、プロセス間の境界にある節点のうち、そのプロセスが担当する節点を○で、そのプロセスが担当しない節点を●で囲みます。まず[1]で、各プロセスは○内の節点の配列Sの値をゼロクリアします(○以外の節点については配列Sに現在の値がそのまま入っています)。

● 次に[2]で図5-2-2の[2]を並列化します(図5-2-7(1)の要素の四隅の矢印参照)。例えば境界上の節点③に着目すると、ランク0では配列Sの現在の値に要素Yの値が加算され(これを③で示します)、ランク1では[1]でゼロクリアしたので配列Sに要素Yの値がそのまま入ります(これを③で示します)。つまりランク0の③にランク1の③を加算すると正しい③になります(これを以下の[3]~[5]で行います)。

● 正しい境界の値を求めるため、[3]~[5]で③を通信して④に加算します(図5-2-7(1)と(2)の間の✓参照)。この部分を図5-2-8の上半分で説明します。図中で[3]などの数字は図5-2-10内の数字に対応します。また例えば③は節点番号(整数)を示し、④は節点③の配列S内の物理量を示します。なお、配列Sは1次元ですが、紙面の都合で2次元のように表します。

まず[3]で、配列Sのうち送信する節点(④内の節点)の物理量を、送信用配列BUF1の宛先プロセスのランクの所に入れます。次に[4]で集団通信ルーチンMPI\_ALLTOALLV(付録参照)を用いて通信を行います。なお、配列ICNT1, ICNT2には各プロセスに送(受)信する個数がセットされていますが、そのプロセスに対して送(受信)を行わない場合はゼロになっている必要があるため、[1]で配列ICNT1, ICNT2をゼロクリアしています。

最後に[5]で、配列BUF2に受信した境界の節点の物理量を、配列Sの該当箇所(④内の節点)に加算します。以上で図5-2-7(2)の④内の節点は正しい物理量になりました。

● [6]で図5-2-2の[3]を並列化します(図5-2-7(2)の◆参照)。**[6]**が終了した時点で、図5-2-7(2)の④の部分には節点の物理量の正しい値が入っていますが、④の部分には正しい値が入っていません。そこで[7]~[9]で④を通信して⑤に代入します(図5-2-7(2)と(3)の間の↘参照)。この部分の通信方法は[3]~[5]と同様なので[15]は加算で、[9]は代入であるところが異なります)、説明は省略します(図5-2-8の下半分参照)。

● [9]が終了した時点で、図5-2-7(3)の④と⑤の両方に節点の物理量の正しい値が入りました。続いて[20]で図5-2-2の[4]を並列化し(図5-2-7(3)の要素の四隅の矢印参照)、[2]で図5-2-2の[5]を並列化します(図5-2-7(4)の◆参照)。

## ■ 計算結果の収集

タイムステップのループが終了し、各プロセスが計算した結果をランク0のプロセスに収集します。

● [20]で図5-2-9(1)に示すように、ランク0以外のプロセスは配列S(節点の物理量)のうち自分が担当した部分を1次元配列BUFにセットします。続いて[23]で、集団通信ルーチンMPI\_GATHERVを用いてランク0の配列BUF1に収集します(付録と、4-6-3-1節の「ケース1」参照)。次に[24]で、ランク0のプロセスは配列BUF1に受信した値を配列Sの該当箇所に代入します。

● [25]では、図5-2-9(2)に示すように、配列Y(要素の物理量)のうち各プロセスが担当した部分をランク0の配列Yに収集します(図で例えば①は要素番号(整数)を示し、[1]は要素①の配列Y内の物理量を示します)。この通信を集団通信ルーチンMPI\_GATHERVで行うと、図5-2-9(2)のランク0の配列Yの空いている部分に集めることができないうため(3-3-5節で説明した、送信バッファと受信バッファがメモリー上で重なってはいけないうという制限により)、MPI\_ISEND, MPI\_IRECVを使用して通信を行いました(4-6-3-1節の「ケース2」参照)。

以上でプログラムの説明を終了します。この方法は修正が面倒ですが、必要最小限のデータのみを通信するため、前節の方法と比較して(一般に)通信量は減少します。なお以下の点に考慮する必要があります。

● 要素番号のつけ方によっては、プロセス間の境界上の節点の数が多くなり、通信量が多くなってしまいう可能性があります。

● 各プロセスが担当する節点数がプロセス間で異なり、ロードバランスが不均等になる可能性があります。

● 図5-2-10では、説明を簡単にするため、配列ISNUM, IBUF1, IBUF2, BUF1, BUF2の大きさが、必要な最小限の大きさよりかなり大きくなっており、使用できるメモリーが少ない場合は改良する必要があります。

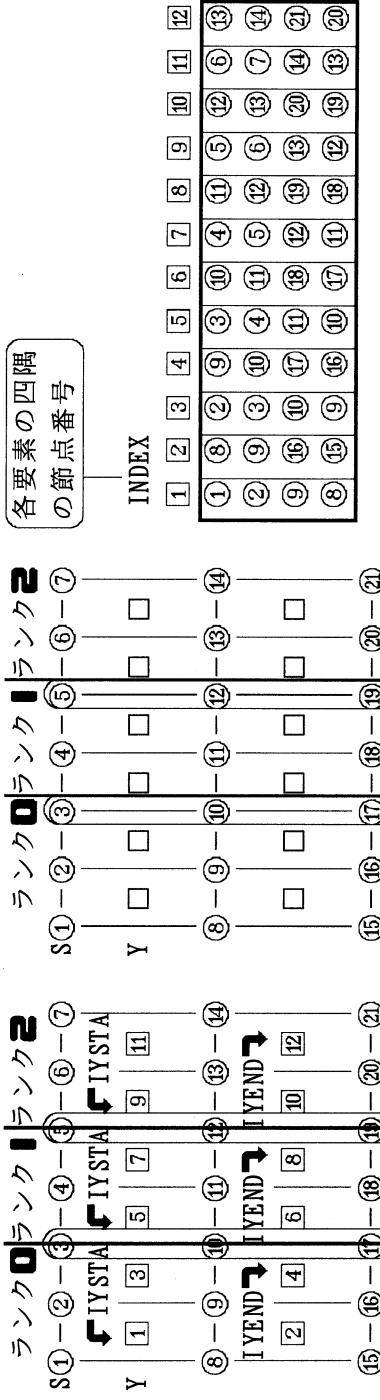


図5-2-6(1) 要素(□)の分割

図5-2-6(2) 節点(○)の分割

BUF1とBUF2内で、各プロセスの  
送受信データが入る最初の位置  
(先頭からの変位)をセットします。

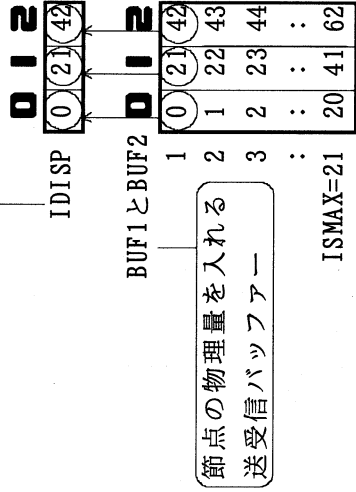
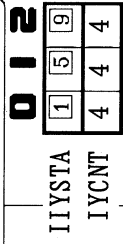


図5-2-6(4)

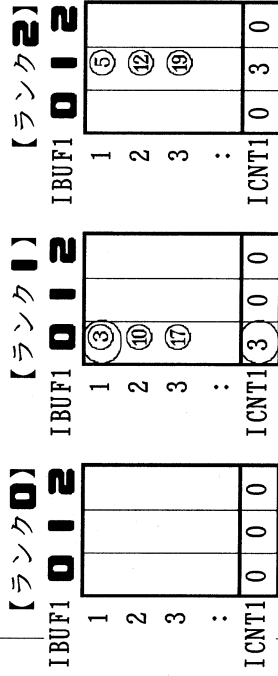
各プロセスが担当  
する最初の要素番号



各プロセスが  
担当する要素数

図5-2-6(5)

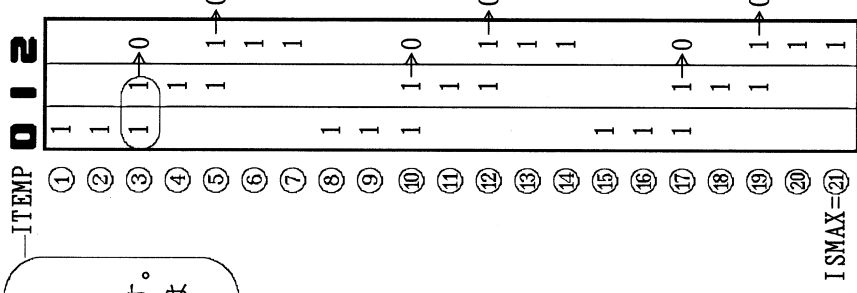
他のプロセスとの境界上の節点のうち、  
自分が担当しない節点の番号が、  
その節点を担当するプロセスの  
ランクの所に入ります。



各プロセスのIBUF1内の節点数

図5-2-6(8)

各プロセスに  
所属する  
節点番号に  
「1」を付けます。  
その他の部分は  
「0」です。

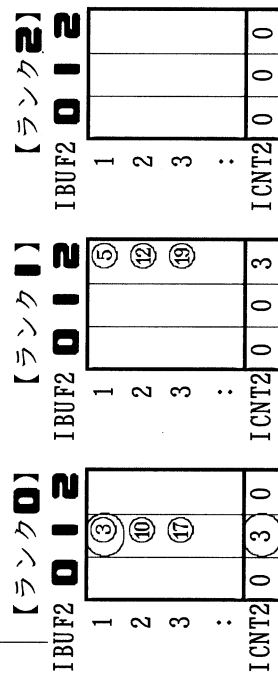


各プロセスが  
担当する節点番号

各プロセスが  
担当する節点数

図5-2-6(7)

他のプロセスとの境界上の節点のうち、  
自分が担当する節点の番号が、  
その節点を担当しないプロセスの  
ランクの所に入ります。



各プロセスのIBUF2内の節点数

図5-2-6(9)

(注)以下の図で

- ○:節点の物理量を入れる配列S
  - □:要素の物理量を入れる配列Y
- を示します。

【ランク0】

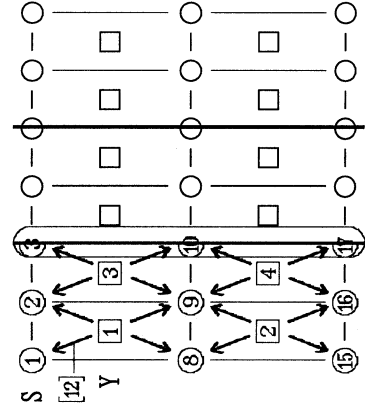


図5-2-7(1)

【ランク1】

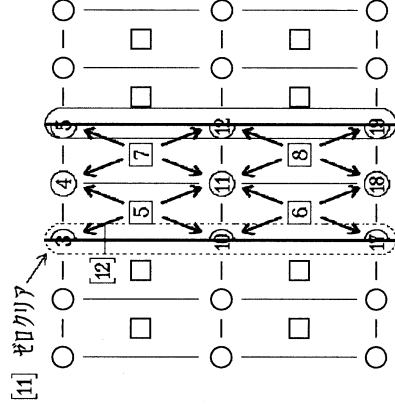


図5-2-7(2)

【ランク2】

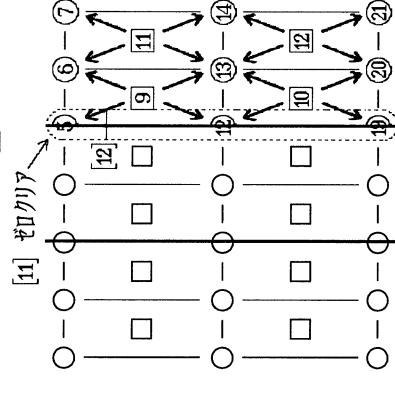


図5-2-7(3)

通信して加算

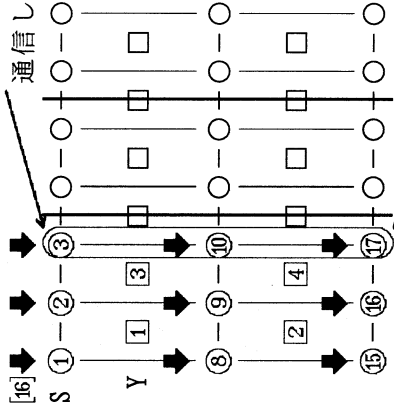


図5-2-7(4)

通信して加算

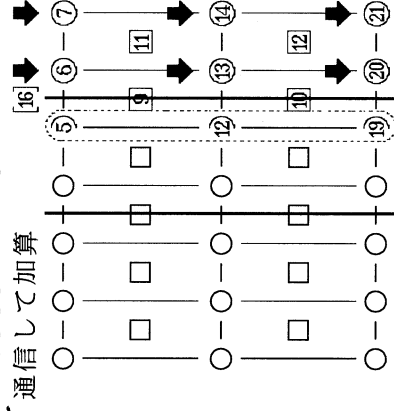


図5-2-7(5)

通信して代入

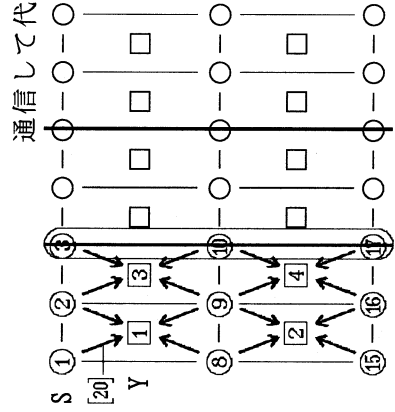


図5-2-7(6)

通信して代入

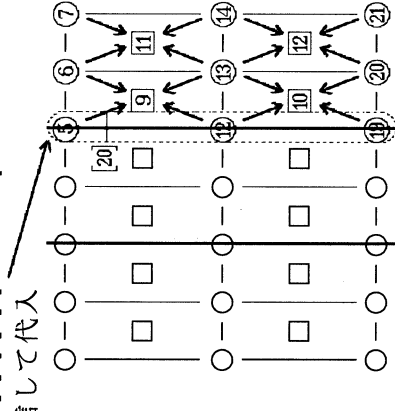


図5-2-7(7)

通信して代入

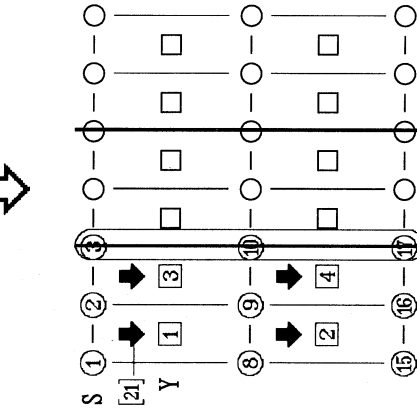


図5-2-7(8)

通信して代入

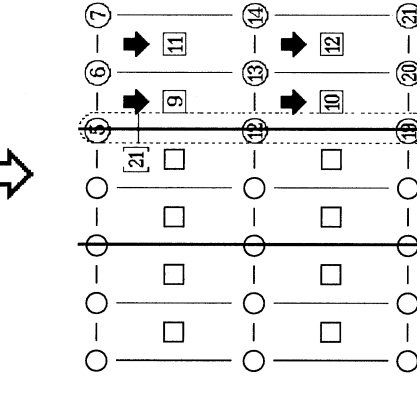


図5-2-7(9)

【ランク0】

S	(1)	(2)	(3)	•	•	•	•
	(8)	(9)	(10)	•	•	•	•
	(15)	(16)	(17)	•	•	•	•

【ランク1】

S	•	•	(3)	(4)	(5)	•	•
	•	•	(10)	(11)	(12)	•	•
	•	•	(17)	(18)	(19)	•	•

【ランク2】

S	•	•	•	•	(5)	(6)	(7)
	•	•	•	•	(12)	(13)	(14)
	•	•	•	•	(19)	(20)	(21)

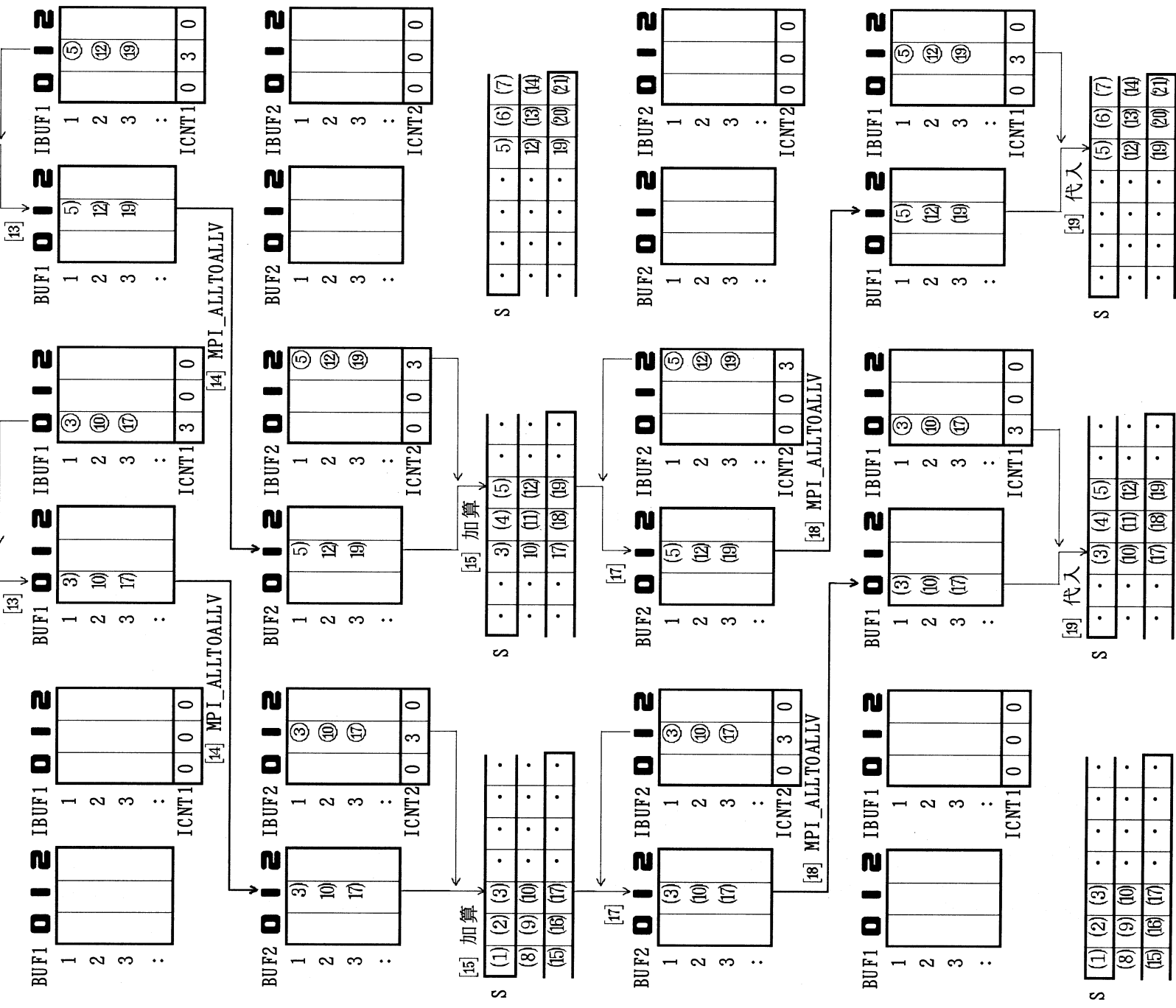


図5-2-8 境界の節点(の物理量)の通信

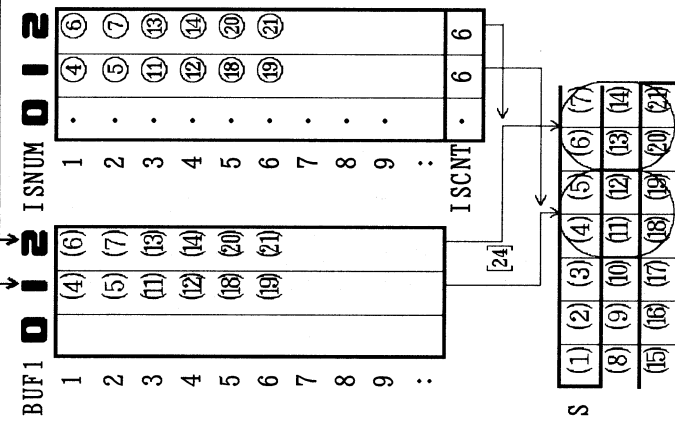
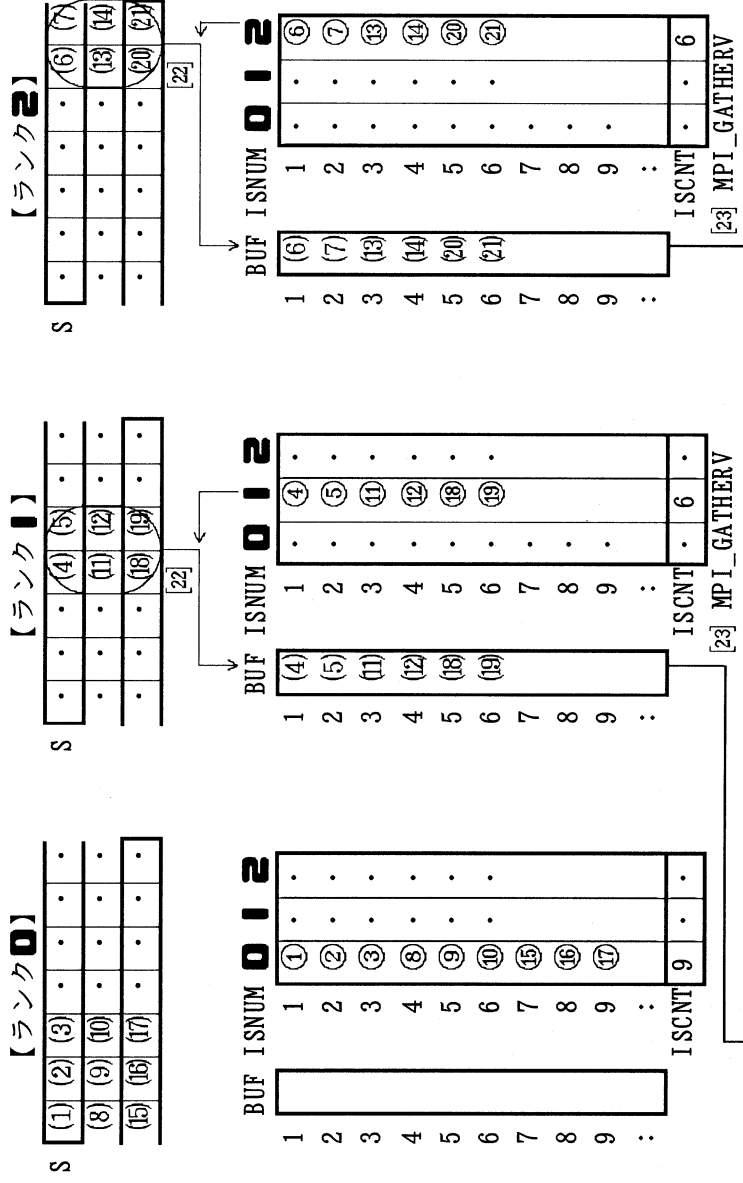


図5-2-9(1) 配列Sの収集

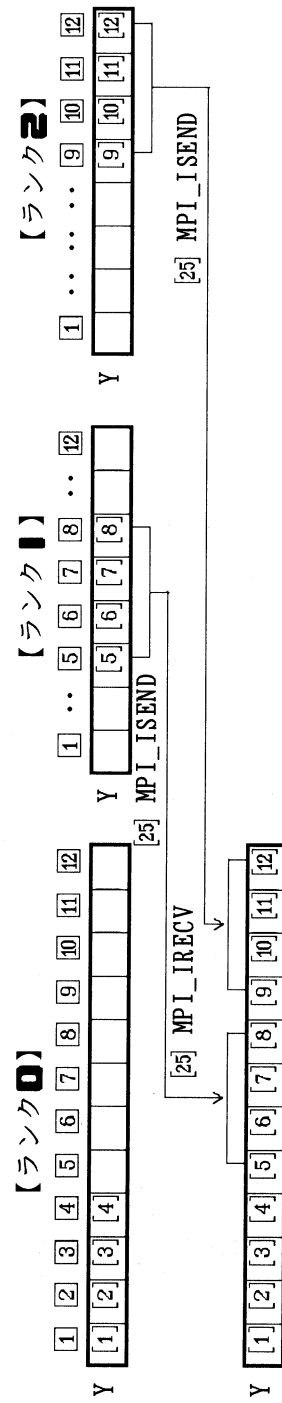


図5-2-9(2) 配列Yの収集



④

```

:
INCLUDE 'mpif.h'
PARAMETER (IYMAX=12, ISMAX=21)
REAL*8 Y(IYMAX), S(ISMAX)
INTEGER INDEX(4, IYMAX)
PARAMETER(NCPU=3)
INTEGER NPROCS, MYRANK
INTEGER ISTATUS(MPI_STATUS_SIZE)
INTEGER ITEMP(ISMAX, 0:NCPU-1)
INTEGER IDISP(0:NCPU-1), ISCNT(0:NCPU-1),
& ISNUM(ISMAX, 0:NCPU-1)
INTEGER IYSTA(0:NCPU-1), IYCNT(0:NCPU-1)
INTEGER ICNT1(0:NCPU-1),
& ICNT2(0:NCPU-1)
INTEGER IBUF1(ISMAX, 0:NCPU-1),
& IBUF2(ISMAX, 0:NCPU-1)
REAL*8 BUF(ISMAX)
REAL*8 BUF1(ISMAX, 0:NCPU-1),
& BUF2(ISMAX, 0:NCPU-1)

```

```

:
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)

```

```

DO IRANK=0, NPROCS-1 → テーブル類の作成
DO I=1, ISMAX
ITEMP(I, IRANK) = 0
ENDDO
ICNT1(IRANK) = 0
ICNT2(IRANK) = 0
ISCNT(IRANK) = 0
ENDDO

```

```

DO IRANK=0, NPROCS-1
IDISP(IRANK) = IRANK*ISMAX
ENDDO
DO IRANK=0, NPROCS-1
CALL PARA_RANGE
& (1, IYMAX, NPROCS, IRANK, IYSTA, IYEND)
IYSTA(IRANK) = IYSTA
IYCNT(IRANK) = IYEND-IYSTA+1
DO IY=IYSTA, IYEND
DO J=1, 4
ITEMP(INDEX(J, IY), IRANK) = 1
ENDDO
ENDDO
ENDDO
CALL PARA_RANGE
& (1, IYMAX, NPROCS, MYRANK, IYSTA, IYEND)

```

④

```

DO I=1, ISMAX
IFLG = 0
DO IRANK=0, NPROCS-1
IF (ITEMP(I, IRANK)==1) THEN
IF (IFLG == 0) THEN
IFLG = 1
IRANK = IRANK
ELSE
ITEMP(I, IRANK) = 0
IF (IRANK == MYRANK) THEN
ICNT1(IRANK)=ICNT1(IRANK)+1 [6]
IBUF1(ICNT1(IRANK), IRANK)=I
ELSEIF (IRANK == MYRANK) THEN
ICNT2(IRANK)=ICNT2(IRANK)+1
IBUF2(ICNT2(IRANK), IRANK)=I
ENDIF
ENDIF
ENDDO
ENDDO

```

```

DO IRANK=0, NPROCS-1
DO I=1, ISMAX
IF (ITEMP(I, IRANK)==1) THEN
ISCNT(IRANK) = ISCNT(IRANK) + 1
ISNUM(ISCNT(IRANK), IRANK) = I
ENDIF
ENDDO
ENDDO

```

```

DO IY=IYSTA, IYEND → 並列計算部分
Y(IY) = IY*10.0
ENDDO
DO I=1, ISCNT(MYRANK)
IS = ISNUM(I, MYRANK)
S(IS) = IS*100.0
ENDDO
DO IY=1, IO
DO IRANK=0, NPROCS-1
DO I=1, ICNT1(IRANK)
S(IBUF1(I, IRANK)) = 0.0
ENDDO
ENDDO
DO IY=IYSTA, IYEND
DO J=1, 4
S(INDEX(J, IY))=S(INDEX(J, IY))+Y(IY) [12]
ENDDO
ENDDO

```

⑤

⑤

```

DO IRANK=0, NPROCS-1
  DO I=1, ICNT1(IRANK)
    BUF1(I, IRANK)=S(IBUF1(I, IRANK)) [13]
  ENDDO
ENDDO
CALL MPI_ALLTOALLV
  (BUF1, ICNT1, IDISP, MPI_REAL8,
  &
  BUF2, ICNT2, IDISP, MPI_REAL8,
  &
  MPI_COMM_WORLD, IERR) [14]
DO IRANK=0, NPROCS-1
  DO I=1, ICNT2(IRANK)
    S(IBUF2(I, IRANK)) =
    S(1BUF2(I, IRANK)) + BUF2(I, IRANK) [15]
  ENDDO
ENDDO
DO I=1, ISCNT(MYRANK)
  IS = ISNUM(I, MYRANK)
  S(IS) = S(IS)*0.25 [16]
ENDDO
DO IRANK=0, NPROCS-1
  DO I=1, ICNT2(IRANK)
    BUF2(I, IRANK)=S(1BUF2(I, IRANK)) [17]
  ENDDO
ENDDO
CALL MPI_ALLTOALLV
  (BUF2, ICNT2, IDISP, MPI_REAL8,
  &
  BUF1, ICNT1, IDISP, MPI_REAL8,
  &
  MPI_COMM_WORLD, IERR) [18]
DO IRANK=0, NPROCS-1
  DO I=1, ICNT1(IRANK)
    S(1BUF1(I, IRANK))=BUF1(I, IRANK) [19]
  ENDDO
ENDDO
DO IY=IYSTA, IYEND
  DO J=1, 4
    Y(IY) = Y(IY) + S(INDEX(J, IY)) [20]
  ENDDO
ENDDO
DO IY=IYSTA, IYEND
  Y(IY) = Y(IY)*0.25 [21]
ENDDO
ENDDO

```

☒5-2-10

⑥

```

IF (MYRANK/=0) THEN
  DO I=1, ISCNT(MYRANK)
    BUF(I) = S(ISNUM(I, MYRANK)) [22]
  ENDDO
ENDIF
CALL MPI_GATHERV
  &
  (BUF, ISCNT(MYRANK), MPI_REAL8,
  &
  BUF1, ISCNT, IDISP, MPI_REAL8,
  &
  0, MPI_COMM_WORLD, IERR) [23]
IF (MYRANK==0) THEN
  DO IRANK=1, NPROCS-1
    DO I=1, ISCNT(IRANK)
      S(ISNUM(I, IRANK)) = BUF1(I, IRANK) [24]
    ENDDO
  ENDDO
ENDIF
IF (MYRANK==0) THEN
  DO IRANK=1, NPROCS-1
    CALL MPI_IRECV(Y(IIYSTA(IRANK)),
    IYCNT(IRANK), MPI_REAL8, IRANK, 1,
    MPI_COMM_WORLD, IREQ, IERR)
    CALL MPI_WAIT(IREQ, ISTATUS, IERR) [25]
  ENDDO
ELSE
  CALL MPI_ISEND(Y(IYSTA), IYEND-IYSTA+1,
  MPI_REAL8, 0, 1, MPI_COMM_WORLD, IREQ, IERR)
  CALL MPI_WAIT(IREQ, ISTATUS, IERR)
ENDIF
IF (MYRANK==0) PRINT *, 'RESULT = ', S, Y
CALL MPI_FINALIZE(IERR)
:
```

⤴

密行列連立一次方程式の直接解法であるLU分解の並列化について説明します。なお、密行列のLU分解のルーチンは並列化が容易なため、6章で紹介する並列版の数値計算ライブラリーに通常含まれており、ライブラリールーチンの方がFortranで作成した並列版よりも一般に高速です。

### ■ 分割方法

ここでは簡単のため $6 \times 6$ の行列で説明します。LU分解では、図5-3-1(1)の左から右に計算ステップが進んでいきます。各ステップでは、図の着色した部分を計算します。このようにLU分解では、ステップが進むにつれて、計算する部分が1行1列ずつ少なくなっていく特徴があります。

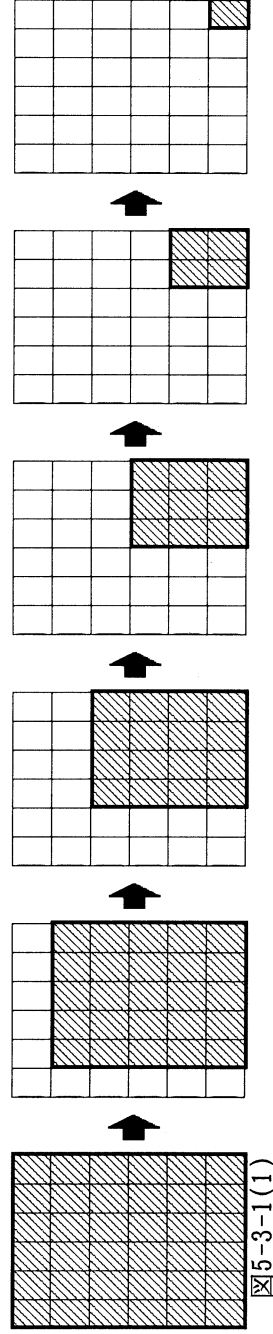


図5-3-1(1)

並列化した場合の分割方法について検討します。ここではプロセス数が3個であるとして、図5-3-1(2)のようにブロック分割で分割した場合、最初のうちはプロセス間のロードバランスは均等ですが、ステップが進むにつれ、ランク0のプロセスの計算部分は終了し、さらに進むとランク1のプロセスの計算部分も終了し、次第にロードバランスが不均等になってしまいます。

一方、サイクリック分割を行った図5-3-1(3)では、どのステップでも、プロセス間のロードバランスは、ほぼ一定になっています。なお、この例では行列の大きさが $6 \times 6$ なのでロードバランスがやや不均等ですが、実際には行列は $1000 \times 1000$ などの大きさなので、ほぼ一定であると言えます。このように、計算が進むにつれて計算部分の大きさが変化する場合、サイクリック分割が有効になります。

なお、並列版の数値計算ライブラリーに含まれるLU分解では、高速化のために行と列の両方向をブロック・サイクリック分割したり、配列を縮小している場合がありますが、以下のプログラムでは説明を分かりやすくするため、列方向のみをサイクリック分割し、配列の縮小はしません。

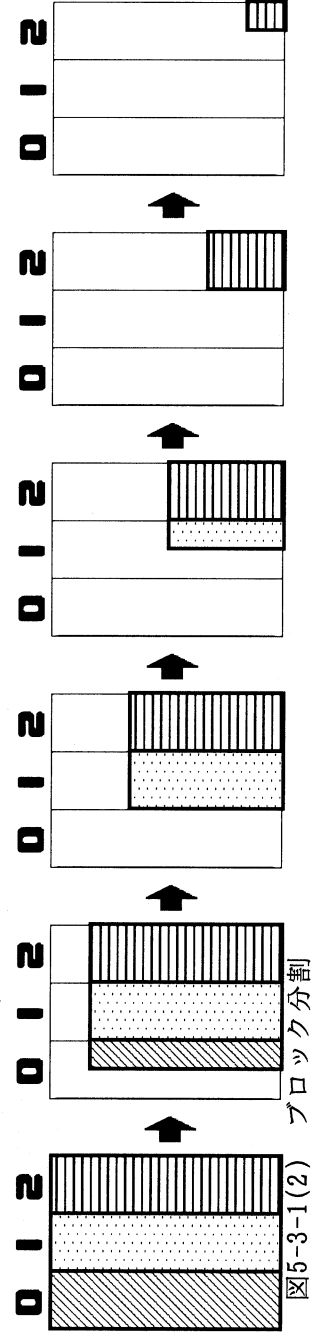


図5-3-1(2) ブロック分割

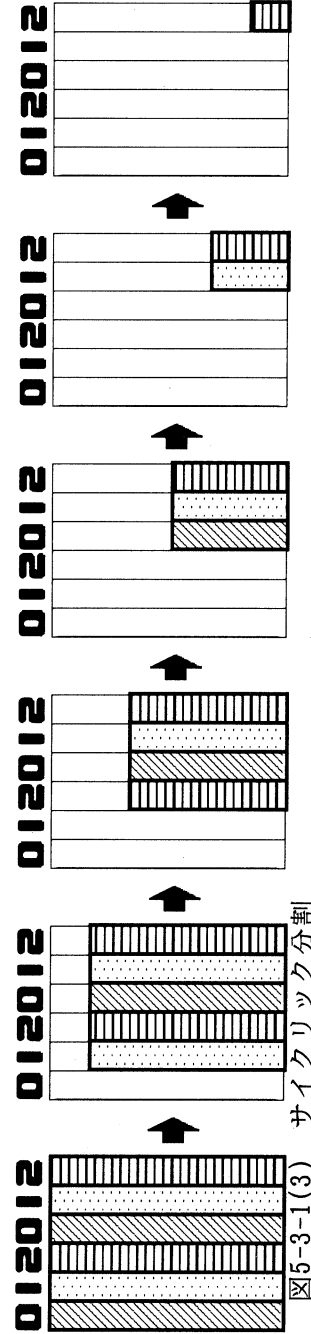


図5-3-1(3) サイクリック分割

## ■ プログラム例

連立一次方程式  $Ax = b$  をLU分解で解く単体版のプログラムを図5-3-2に示します。なお、説明を分かりやすくするため、ループ・アンローリング、ブロック化などのチューニング、およびピボティングは行っていません。

LU分解では、以下に示すように、まず連立一次方程式  $Ax = b$  の行列  $A$  を、 $L$   $U$  という2つの三角行列に分解します。次に行列  $L$   $U$  とベクトル  $b$  を使って求解(前進消去と後退代入)を行い、作業ベクトル  $y$  を経て答えのベクトル  $x$  が求まります。

```
A = L U      【分解】
y = L-1b    【求解(前進消去)】
x = U-1y    【求解(後退代入)】
```

上記の行列  $A$  と  $L$   $U$  が図5-3-2の配列  $A$  に対応し、ベクトル  $b$ ,  $y$ ,  $x$  が配列  $B$  に対応します。つまり最初設定されていた行列  $A$  の値は壊されて  $L$   $U$  (行列の  $\blacksquare$  部分が  $L$ 、 $\blacktriangle$  部分が  $U$ ) になります。ベクトル  $b$  の値も壊されて作業ベクトル  $y$  になり、最後に答えのベクトル  $x$  になります。

求解部分の解法には、図5-3-2に示す内積型と、図5-3-3に示す三項演算型があります。内積型はメモリー上で飛び飛びの要素を参照するため、スカラ計算機ではキャッシュミスが発生して速度が低下しますが(参考文献[6]の4章参照)、三項演算型より並列化がしやすいので、ここでは内積型を使用します。分解より求解の方が計算時間ははるかに短いので、分解、求解をともに一度だけ行う場合は内積型で問題ないですが、分解が1回のみで求解だけを何度も行う場合は、内積型は速度が遅いので使用しないで下さい。

```
PROGRAM MAIN
PARAMETER(N=6)
DIMENSION A(N,N),B(N)
:
① DO K=1,N      【分解】
②   DO I=K+1,N
      A(I,K) = A(I,K)/A(K,K)
      ENDDO
③   DO J=K+1,N
      DO I=K+1,N
        A(I,J) = A(I,J) - A(I,K)*A(K,J)
      ENDDO
      ENDDO
④ DO I=1,N      【求解(前進消去)】
      SUM = 0.0
      DO J=1,I-1
        SUM = SUM + A(I,J)*B(J)
      ENDDO
      B(I) = B(I) - SUM
      ENDDO
⑤ DO I=N,1,-1  【求解(後退代入)】
      SUM = 0.0
      DO J=I+1,N
        SUM = SUM + A(I,J)*B(J)
      ENDDO
      B(I) = (B(I)-SUM)/A(I,I)
      ENDDO
:
```

図5-3-2 求解が内積型(単体版)

```
:
⑥ DO J=1,N      【求解(前進消去)】
      DO I=J+1,N
        B(I) = B(I) - A(I,J)*B(J)
      ENDDO
      ENDDO
⑦ DO J=N,1,-1  【求解(後退代入)】
      B(J) = B(J)/A(J,J)
      DO I=1,J-1
        B(I) = B(I) - A(I,J)*B(J)
      ENDDO
      ENDDO
:
```

図5-3-3 求解が三項演算型(単体版)

並列化したプログラムを以下に示し、次のページ以降で説明します。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
PARAMETER(N=6)
DIMENSION A(N,N),B(N)
INTEGER MAP(N)

CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD,NPROCS,IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD,MYRANK,IERR)
DO I=1,N
MAP(I) = MOD(I-1,NPROCS)
ENDDO
:
DO K=1,N 【分解】
IF (MAP(K)==MYRANK) THEN
DO I=K+1,N
A(I,K) = A(I,K)/A(K,K)
ENDDO
ENDIF
CALL MPI_BCAST(A(K+1,K),N-K,MPI_REAL,
MAP(K),MPI_COMM_WORLD,IERR)
&
DO J=K+1,N
IF (MAP(J) == MYRANK) THEN
DO I=K+1,N
A(I,J) = A(I,J) - A(I,K)*A(K,J)
ENDDO
ENDIF
ENDDO
ENDDO

```

```

DO I=1,N 【求解(前進消去)】
SUM = 0.0
DO J=1,I-1
IF (MAP(J)==MYRANK)
SUM = SUM + A(I,J)*B(J)
ENDDO
CALL MPI_REDUCE(SUM,SSUM,1,MPI_REAL,
MPI_SUM,MAP(I),MPI_COMM_WORLD,IERR)
&
IF (MAP(I)==MYRANK) B(I) = B(I) - SSUM
ENDDO
DO I=N,1,-1 【求解(後退代入)】
SUM = 0.0
DO J=I+1,N
IF (MAP(J)==MYRANK)
SUM = SUM + A(I,J)*B(J)
ENDDO
CALL MPI_REDUCE(SUM,SSUM,1,MPI_REAL,
MPI_SUM,MAP(I),MPI_COMM_WORLD,IERR)
&
IF (MAP(I)==MYRANK)
B(I) = (B(I)-SSUM)/A(I,I)
ENDDO
DO I=1,N
IF (MAP(I)/=MYRANK) B(I) = 0.0
ENDDO
CALL MPI_ALLREDUCE(MPI_IN_PLACE,B,N,
MPI_REAL,MPI_SUM,MPI_COMM_WORLD,IERR)
:

```

図5-3-4 求解が内積型(並列版)

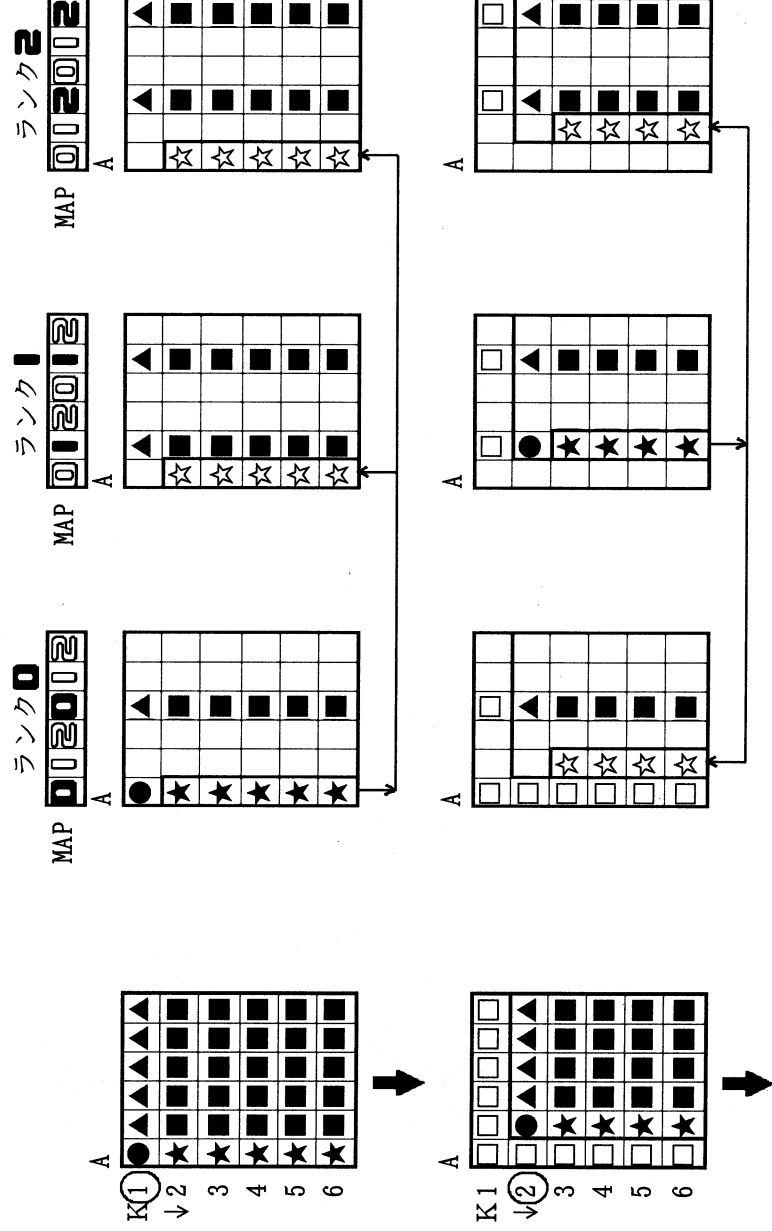
## ■ 分解

### 【単体版の動作】

- 図5-3-5(1)に、6×6の行列Aを行列L Uに変換する過程の一部を示します。
- 図5-3-2の①のループでK=1のとき、②で、図5-3-5(1)(上段)の各★に対し、「★ = ★/●」を行います。
- 次に③で、図5-3-5(1)(上段)の各■に対し、「■ = ■ - ★\*▲」を行います。
- ①のループでK=2となり、図5-3-5(1)(下段)の太線内に対し、上記と同様の処理を行います。
- 以後①のループでK=3,4,5と同様の処理を行います。前述のように、計算する部分は次第に少なくなります。

### 【並列版の動作】

- 前述のように、列方向をサイクリック分割して並列化を行います。
- 図5-3-4の①で、行列Aの各列番号と、その列を担当するプロセスのランク値の対応を表す、MAPという対応表(図5-3-5(2)(上段))を作成します(4-5-5節の分割方法2参照)。なお、MAPを使用すると後述するようにIF文が必要となり、計算部分の速度が低下します。MAPを使用しない方法(4-5-5節の分割方法1)もあります。ここでは説明を分かりやすくするために使用しました。
- ②のループでK=1のとき、④で、図5-3-6(2)(上段)の1列目を担当するランク0のみが、各★に対し、「★ = ★/●」を行います。③のIF文があるため、ランク0のみがこの処理を行います。●⑤で、ランク0からその他のプロセスに、計算した★のデータを送信します。このとき送信元のプロセスは、1列目を担当するMAP(1)(ランク0)のプロセスとなります。
- ⑦で、各プロセスは自分が担当する列の各■に対し、「■ = ■ - ★\*▲」を行います(ランク0以外のプロセスでは「■ = ■ - ☆\*▲」)。⑥のIF文があるため、各プロセスは自分が担当する列のみを処理します。
- ②のループでK=2となり、図5-3-5(2)(下段)の太線内に対し、上記と同様の処理を行います。⑤の通信で、送信元のプロセスは、2列目を担当するMAP(2)(ランク1)のプロセスとなります。
- 以後②のループでK=3,4,5と同様の処理を行います。⑤の通信で、送信元のプロセスはランク2,0,1と変化します。



(以下略)

図5-3-5(1) 単体版の動作

図5-3-5(2) 並列版の動作

■ 求解(前進消去)

【単体版の動作】

- 図5-3-6(1)に、ベクトルbを作業ベクトルyに変換する過程の一部を示します。
- 説明の都合上、図5-3-2の④のループで、I=1,2,3,4の順に以下と同様の処理がすでに行われ、B(1)~B(4)は計算が終了しているとしします。
- ④のループでI=5のとき、⑤で、図5-3-6(1)(上段)の□と○の内積を計算し、結果が変数SUMに入ります。
- ⑥で★に対し、「★ = ★ - SUM」を行います。これでB(5)の計算が終了しました。
- ④のループでI=6となり、図5-3-6(1)(下段)に示すように上記と同様の処理をし、B(6)を計算します。

【並列版の動作】

- LU分解の結果、図5-3-6(2)(上段)に示すように、配列Aのうち、対応表MAPに従って各プロセスが担当した列のみに正しい値(□と■)が入っています。配列Bの各要素も、対応表MAPに従って上から順にランク○,■,2,0,1,2のプロセスが担当します。
- 単体版の説明と同様に、図5-3-4の⑧のループでI=5の場合から説明します。⑨で、各プロセスは図5-3-6(2)(上段)の□と○のうち、自分が担当する要素の内積を計算し、結果が変数SUMに入ります。⑩のIF文があるため、各プロセスは自分が担当する要素のみを計算します。前述のように、IF文を使用すると速度が低下しますが、ここでは説明を分かりやすくするため使用しました。
- 各プロセスが求めた部分的な内積SUMの値を、⑭で通信しながら合計し、結果が★を担当するランク■のプロセスの変数SSUMに送られます。
- ⑮で、ランク■のプロセスは★に対し、「★ = ★ - SSUM」を行います。これでB(5)の計算が終了しました。
- ⑧のループでI=6となり、図5-3-6(2)(下段)に示すように上記と同様の処理をし、B(6)を計算します。⑯の通信で、宛先のプロセスは★を担当するランク2となります。

(以上略)

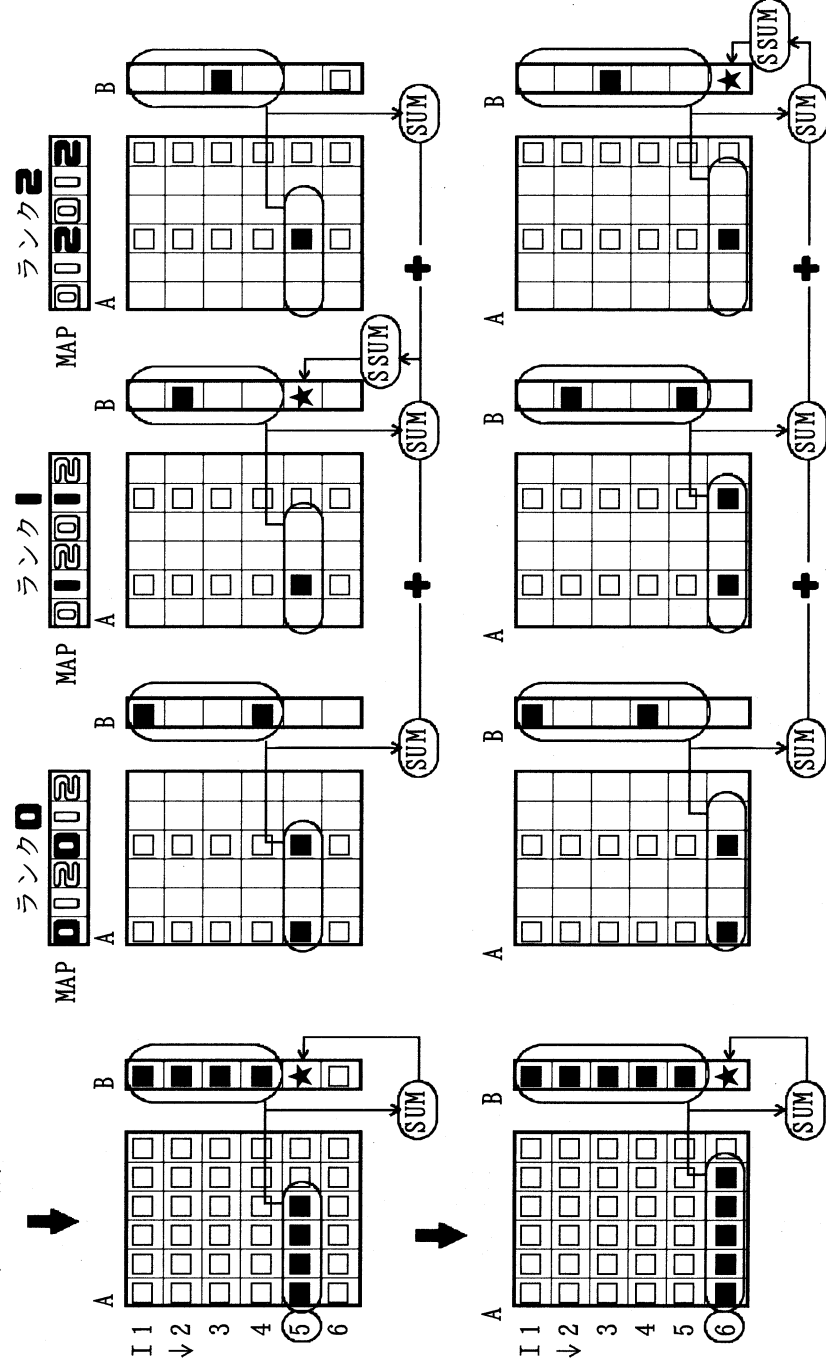


図5-3-6(1) 単体版の動作

図5-3-6(2) 並列版の動作

■ 求解(後退代入)

【単体版の動作】

- 図5-3-7(1)に、作業ベクトル $\gamma$ を答えのベクトル $x$ に変換する過程の一部を示します。
- 説明の都合上、図5-3-2の⑦のループで、 $I=6, 5, 4, 3$ の順に以下と同様の処理がすで行われ、 $B(6) \sim B(3)$ は計算が終了してしているとします。
- ⑦のループで $I=2$ のとき、⑧で、図5-3-7(1)(上段)の○と□の内積を計算し、結果が変数SUMに入ります。
- ⑨で★に対し、「★ = (★-SUM)/●」を行います。これでB(2)の計算が終了しました。
- ⑦のループで $I=1$ となり、図5-3-7(1)(下段)に示すように上記と同様の処理をし、B(1)を計算します。

【並列版の動作】

- 単体版の説明と同様に、図5-3-4の④のループで $I=2$ の場合から説明します。⑬で、各プロセスは図5-3-7(2)(上段)の○と□のうち、自分が担当する要素の内積を計算し、結果が変数SUMに入ります。⑭のIF文があるため、各プロセスは自分が担当する要素のみを計算します。前述のように、IF文を使用すると速度が低下しますが、ここでは説明を分かりやすくするため使用しました。
- 各プロセスが求めた部分的な内積SUMの値を、⑮で通信しながら合計し、結果が★を担当するランク■のプロセスの変数SSUMに送られます。
- ⑮で、ランク■のプロセスは★に対し、「★ = (★-SSUM)/●」を行います。これでB(2)の計算が終了しました。
- ⑬のループで $I=1$ となり、図5-3-7(2)(下段)に示すように上記と同様の処理をし、B(1)を計算します。⑯の通信で、宛先のプロセスは★を担当するランク□となります。
- 後退代入が終了し、最後に各プロセスが担当した配列Bの要素を1箇所に集めます。⑰で、各プロセスは配列Bのうち自分が担当しない要素に0を入れます。そして⑱で、図5-3-7(2)(下段)に示すように、配列Bの各要素を通信しながら合計し(4-6-6節参照)、結果が全プロセスの配列Bに送られます。これが答えのベクトル $x$ となります。

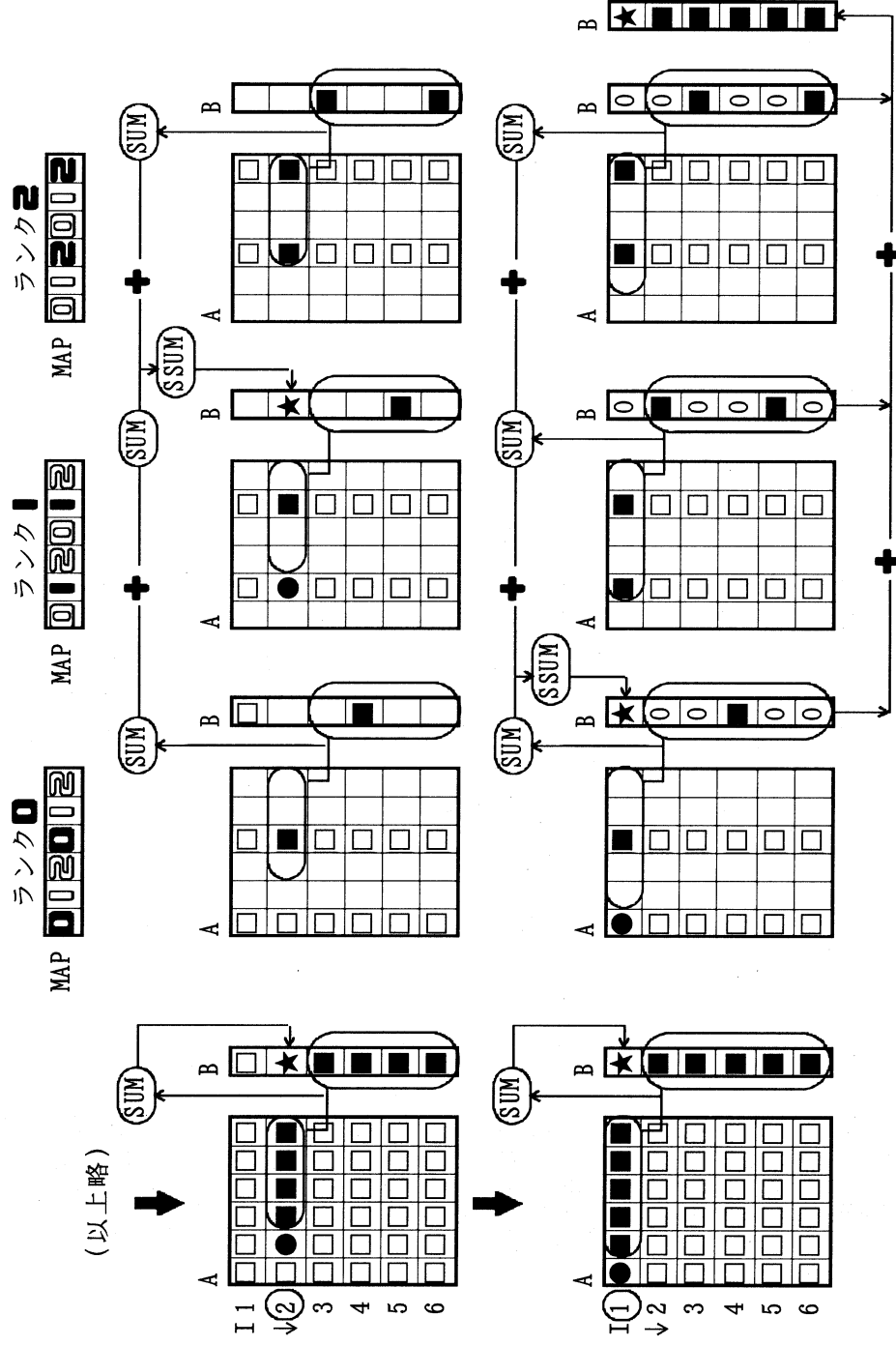


図5-3-7(1) 単体版の動作

図5-3-7(2) 並列版の動作



## 5-4 ICCG法

本節では、差分法などで現れる対称帯行列の反復解法であるICCG法を並列化する方法として、パイプライン法とパラレル・ブロック・オーダリング法を紹介します。これらの方法は非対称帯行列の反復解法(ILUBCG法、ILUCGS法、ILU Bi-CGSTAB法など)にも適用することができます。

なお、第6章で紹介する並列版の数値計算ライブラリーにも適用することができます。(例えばPCP)、実用上はそちらを使用した方がよいと思われま

す。本書では紹介しませんが、ICCG法の前処理部分を簡単にしたSCG法(参考文献[49])は、前処理が簡単なで収束性にやや難がありますが、並列計算に向いています。さらに、ベクトル化効率を高めるためのTF法(参考文献[3])は、並列化も可能ではないかと思われま

### 5-4-1 パイプライン法による並列化

#### ■ ICCG法のアルゴリズム

まずICCG法をパイプライン法(4-6-7節参照)で並列化する方法を説明します。単体版のICCG法のアルゴリズムを図5-4-1(1)に、プログラムを図5-4-1(2)に示します。これらは参考文献[22][23]を参照しました。なお、参考文献で説明されているMICCG法(ICCG法の加速版)のロジックは、説明を簡単にするために省略していますので、図5-4-1(2)のプログラムはあくまでサンプルと考えると下さい。

図5-4-1(1)で行列A、ベクトルx、bは連立一次方程式Ax=bを表し、ベクトルp,q,rは作業用ベクトル、 $\alpha, \beta, c_1, c_2, c_3$ はスカラー変数です。実線、二重線、波線の計算はサブルーチンを使用します。

```

① x0, p0, q, r0を初期化します。
② A ≐ L D LT (サブルーチンDECOMPを使用)
③ q = A x0 (サブルーチンAXSUBを使用)
④ r0 = b - q
⑤ p0 = (L D LT)-1 r0 (サブルーチンLDLTを使用)
⑥ c1 = (r0, p0)
D0 k = 0, 1, ... (収束反復のループ)
⑦ q = A pk (サブルーチンAXSUBを使用)
⑧ c2 = (pk, q)
⑨  $\alpha_k = c_1 / c_2$ 
⑩ xk+1 = xk +  $\alpha_k$  pk
⑪ rk+1 = rk -  $\alpha_k$  q
⑫  $\|x_{k+1} - x_k\| / \|x_k\| < \epsilon$  のとき終了します。
⑬ q = (L D LT)-1 rk+1 (サブルーチンLDLTを使用)
⑭ c3 = (rk+1, q)
⑮  $\beta_k = c_3 / c_1$ 
⑯ c1 = c3
⑰ pk+1 = q +  $\beta_k$  pk
ENDDO

```

図5-4-1(1)

■ 単体版プログラムの説明

図5-4-1(2)のプログラムを説明します。なお、図5-4-1(1)内の番号と図5-4-1(2)内の番号は対応しています。図5-4-1(2)の[1]で連立一次方程式  $Ax = b$  の  $A$  と  $b$  に値を設定し(設定方法については後述します)、[2]で計算条件を設定し、[3]でICCG法の計算を実際に行うサブルーチンICCGをコールします。  
サブルーチンICCGの各引数の意味を以下に示します(配列の構造は後述します)。

- A, X, B : 連立一次方程式  $Ax = b$  に相当します。AとBを設定して本サブルーチンをコールすると、解がXに戻ります。
- D : 図5-4-1(1)で  $A \div LDL^T$  と分解したときの対角要素Dに相当する配列です。
- MX, MY : 図5-4-2(4)に示すように、行列Aのバンド幅がMX、行列Aの行数がMX×MYとなります。
- P, Q, R : 図5-4-1(1)の作業用ベクトルp, q, rに相当する配列です。
- EPS : 収束判定値を設定します。サブルーチンが終了すると、収束したときの相対誤差  $\|x_{k+1} - x_k\| / \|x_k\|$  が戻ります。
- ITR : 打ち切りまでの最大繰り返し回数を設定します。サブルーチンが終了すると、収束したときの繰り返し回数に戻ります。
- IER : サブルーチンが終了すると、エラーコードが戻ります。IER=0ならば正常終了、IER=1ならば収束しなかったことを意味します。

<pre> PROGRAM MAIN IMPLICIT REAL*8(A-H,0-Z) PARAMETER (MX=6,MY=6) DIMENSION A(0:MX+1,0:MY+1,3),B(MX,MY) DIMENSION X(0:MX+1,0:MY+1) DIMENSION D(0:MX+1,0:MY+1),P(0:MX+1,0:MY+1) DIMENSION Q(0:MX+1,0:MY+1),R(0:MX+1,0:MY+1) : 配列AとBに値を設定します。 [1] EPS=1.0D-7 [2] ITR=1000 [2] CALL ICCG(A, MX, MY, B, EPS, ITR, X, D, P, Q, R, IER) [3] : END </pre>	<div style="text-align: center;">④ →</div> <pre> SUBROUTINE ICCG(A, MX, MY, B, EPS, ITR, &amp; X, D, P, Q, R, IER) IMPLICIT REAL*8(A-H,0-Z) DIMENSION A(0:MX+1,0:MY+1,3), B(MX,MY) DIMENSION X(0:MX+1,0:MY+1) DIMENSION D(0:MX+1,0:MY+1), P(0:MX+1,0:MY+1) DIMENSION Q(0:MX+1,0:MY+1), R(0:MX+1,0:MY+1) DO J = 0, MY+1 DO I = 0, MX+1 D(I,J) = 0.0D0 X(I,J) = 0.0D0 P(I,J) = 0.0D0 Q(I,J) = 0.0D0 R(I,J) = 0.0D0 ENDDO ENDDO CALL DECOMP(A,D, MX, MY) CALL AXSUB(A,X,Q, MX, MY) DO J = 1, MY DO I = 1, MX R(I,J) = B(I,J) - Q(I,J) ENDDO ENDDO CALL LDLT(A,D,R,P, MX, MY) C1 = 0.0D0 DO J = 1, MY DO I = 1, MX C1 = C1 + R(I,J)*P(I,J) ENDDO ENDDO </pre> <div style="text-align: right;">← ③</div>
--	---



収束反復ループ

```

D0 K = 1, ITR
CALL AXSUB(A,P,Q,MX,MY)
C2 = 0.0D0
D0 J = 1, MY
  D0 I = 1, MX
  C2 = C2 + P(I,J)*Q(I,J)
  ENDDO
ENDDO
ALPHA = C1 / C2
X1 = 0.0D0
X2 = 0.0D0
D0 J = 1, MY
  D0 I = 1, MX
  Y = X(I,J)
  X(I,J) = X(I,J) + ALPHA*P(I,J)
  R(I,J) = R(I,J) - ALPHA*Q(I,J)
  X1 = X1 + Y*Y
  X2 = X2 + (X(I,J)-Y)**2
  ENDDO
ENDDO
IF (X1 / = 0.0) THEN
  RES = DSQRT(X2/X1)
  IF (RES <= EPS) THEN
    ITR = K
    EPS = RES
    IER = 0
    RETURN
  ENDF
ENDF
CALL LDLT(A,D,R,Q,MX,MY)
C3 = 0.0D0
D0 J = 1, MY
  D0 I = 1, MX
  C3 = C3 + R(I,J)*Q(I,J)
  ENDDO
ENDDO
BETA = C3/C1
C1 = C3
D0 J = 1, MY
  D0 I = 1, MX
  P(I,J) = Q(I,J) + BETA*P(I,J)
  ENDDO
ENDDO
IER = 1
EPS = RES
END
  
```



図5-4-1(2)



```

SUBROUTINE DECOMP(A,D,MX,MY)
IMPLICIT REAL*8(A-H,0-Z)
DIMENSION A(0:MX+1,0:MY+1,3)
DIMENSION D(0:MX+1,0:MY+1)
D0 J = 1, MY
  D0 I = 1, MX
  D(I,J) = 1.0D0/( A(I,J,3)
    & -D(I,J-1)*A(I,J,1)**2
    & -D(I-1,J)*A(I,J,2)**2 )
  ENDDO
ENDDO
END

SUBROUTINE AXSUB(A,X,Y,MX,MY)
IMPLICIT REAL*8(A-H,0-Z)
DIMENSION A(0:MX+1,0:MY+1,3)
DIMENSION X(0:MX+1,0:MY+1),Y(0:MX+1,0:MY+1)
D0 J = 1, MY
  D0 I = 1, MX
  Y(I,J) = A(I,J,1)*X(I,J-1)
    & + A(I,J,2)*X(I-1,J)
    & + A(I,J,3)*X(I,J)
    & + A(I+1,J,2)*X(I+1,J)
    & + A(I,J+1,1)*X(I,J+1)
  ENDDO
ENDDO
END

SUBROUTINE LDLT(A,D,X,Y,MX,MY)
IMPLICIT REAL*8(A-H,0-Z)
DIMENSION A(0:MX+1,0:MY+1,3)
DIMENSION D(0:MX+1,0:MY+1)
DIMENSION X(0:MX+1,0:MY+1),Y(0:MX+1,0:MY+1)
D0 J = 1, MY
  D0 I = 1, MX
  Y(I,J) = D(I,J)*( X(I,J)
    & -A(I,J,1)*Y(I,J-1)
    & -A(I,J,2)*Y(I-1,J) )
  ENDDO
ENDDO
D0 J = MY, 1, -1
  D0 I = MX, 1, -1
  Y(I,J) = Y(I,J)-D(I,J)*(
    & A(I+1,J,2)*Y(I+1,J)
    & +A(I,J+1,1)*Y(I,J+1) )
  ENDDO
ENDDO
END
  
```

行列Aとして、本節では2次元差分法で現れる5重対角対称帯行列を想定します。例えば図5-4-2(1)のような計算領域の場合、行列Aは図5-4-2(4)のように(対称なので)3本の帯 $a_i, b_i, c_i$ から構成されます(紛らわしいですが、 $a_i, b_i, c_i$ は図5-4-1(1)の $b$ や $c$ とは異なることに注意して下さい)。

通常ICCG法のプログラムでは $a_i, b_i, c_i$ や図5-4-1(1)の $b, x, p, q, r$ を1次元ベクトルとして保持しますが、並列化する場合は、2次元配列(つまり図5-4-2(1)の計算領域と同じ形状)で扱った方が処理が簡単なので、以後これらのベクトルを2次元配列で表現することにします。このとき図5-4-2(2)のように要素の添字は1次元ベクトルのまま $a_1, a_2, \dots$ と表すこともできますが、2次元配列として $a_{11}, a_{12}, \dots$ のように表した方が自然なので、図5-4-2(3)のように表すことにします。

プログラム内で使用している実際の配列名と構造を図5-4-3に示します。各配列(Bを除く)は計算領域よりも周囲が1要素分大きくなくっていることに注意して下さい(理由は各サブルーチンでの計算を簡単にするためです)。図5-4-1(2)の[1]では、配列AとBに値を設定しますが、図5-4-3の配列A内で『0』となっている部分には(配列の周囲の部分も含め)必ずゼロを入れて下さい。なお、図5-4-1(1)のその他のベクトルD, p, q, rは図5-4-3の配列Xと同じ構造をしています。

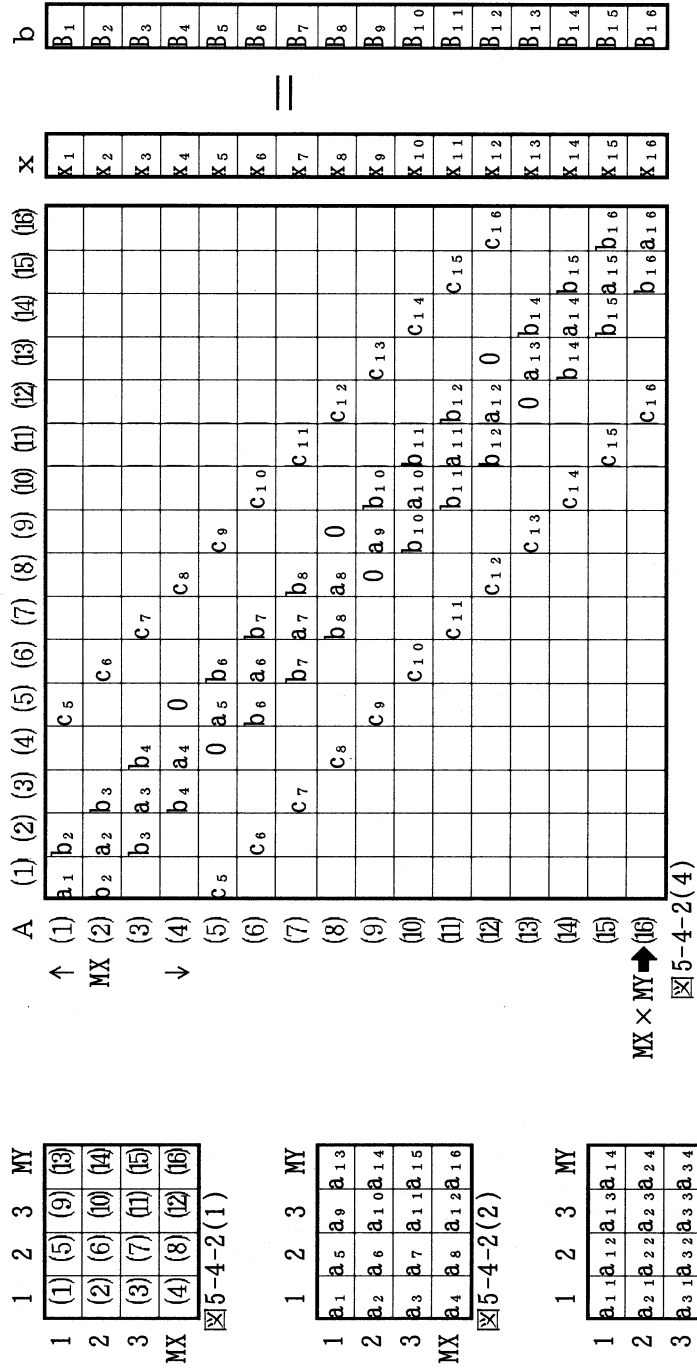


図5-4-2(4)

A(I, J, 1) (=c<sub>i</sub>)

0	1	2	3	MY	MY+1
0	0	0	0	0	0
1	0	0	$c_{12}$	$c_{13}$	$c_{14}$
2	0	0	$c_{22}$	$c_{23}$	$c_{24}$
3	0	0	$c_{32}$	$c_{33}$	$c_{34}$
MX	0	0	$c_{42}$	$c_{43}$	$c_{44}$
MX+1	0	0	0	0	0

A(I, J, 2) (=b<sub>i</sub>)

0	1	2	3	MY	MY+1
0	0	0	0	0	0
1	0	0	0	0	0
2	0	$b_{22}$	$b_{23}$	$b_{24}$	0
3	0	$b_{32}$	$b_{33}$	$b_{34}$	0
MX	0	$b_{42}$	$b_{43}$	$b_{44}$	0
MX+1	0	0	0	0	0

A(I, J, 3) (=a<sub>i</sub>)

0	1	2	3	MY	MY+1
0	0	0	0	0	0
1	0	$a_{11}$	$a_{12}$	$a_{13}$	$a_{14}$
2	0	$a_{21}$	$a_{22}$	$a_{23}$	$a_{24}$
3	0	$a_{31}$	$a_{32}$	$a_{33}$	$a_{34}$
MX	0	$a_{41}$	$a_{42}$	$a_{43}$	$a_{44}$
MX+1	0	0	0	0	0

X(I, J) (=x, D, p, q, r)

0	1	2	3	MY	MY+1
0	0	0	0	0	0
1	0	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$
2	0	$x_{21}$	$x_{22}$	$x_{23}$	$x_{24}$
3	0	$x_{31}$	$x_{32}$	$x_{33}$	$x_{34}$
MX	0	$x_{41}$	$x_{42}$	$x_{43}$	$x_{44}$
MX+1	0	0	0	0	0

図5-4-3



■ 配列の分割方法

以後、サブルーチンICCG内の各部の並列化について説明します。並列化後のプログラムを図5-4-10に示します(図5-4-1(1)内の番号と図5-4-10内の番号は対応しています)。以後、計算領域を6×6とし、図5-4-3の各配列を図5-4-4のように表現します。

まず配列の分割方法ですが、パイプライン法では1次元目または2次元目のどちらから一方でブロック分割します。本節では説明を簡単にするため配列の縮小は行っていません。このため通信の容易さから2次元目でブロック分割することになります。各プロセスの担当部分を図5-4-4に示します(図中のJSTA, JENDはランクのプロセスの範囲を示します)。この分割を図5-4-10の[5]で行います。

実際のプログラムでは、配列の大きさが長方形の場合、境界の長さが短くなる(通信量が少なくなる)次元で分割するのがよいでしょう(4-5-8節参照)。

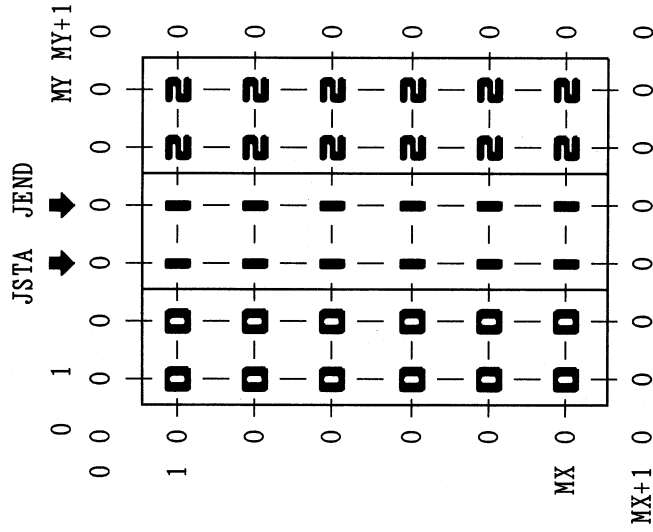


図5-4-4

■ ベクトル同志の演算と内積計算の並列化

図5-4-1(1)のうち、①,④,⑥,⑧,⑩,⑫,⑭,⑯はベクトル同志の演算や内積計算なので、容易に並列化することができます。このうち①以外の部分を並列化しました(図5-4-10の該当番号参照)。それに伴い⑥,⑧,⑩,⑫,⑭の内積(や合計の)計算の後、求めた部分合和を[6]~[9]で通信しながら合計し、結果を全プロセスに渡します。なお、3-3-6節で説明したように、内積や合計の計算部分を並列化すると計算結果が若干変わります。

計算が収束したら最後に、各プロセスは、自分が求めた解ベクトルxの値を[10]で他の全プロセスに送信します。あらかじめ[4]で、各プロセスの担当部分の派生データ型を、3-5-3-2節で説明した自作のサブルーチンPARAMETER\_BLOCK2で作成して配列ITYPEに代入し(MPI-2が使用できるマシン環境の場合は、MPI-2で提供されているMPI\_TYPE\_CREATE\_SUBARRAY(3-5-3-1節参照)を使用して下さい)、4-6-3-2節の方法で全プロセスに通信します。なお、[3]以降の処理で解ベクトルxの値を例えばランク0のプロセスしか使用しないのであれば、4-6-3-1節の方法でランク0に収集します。

## ■ サブルーチン D E C O M P

本サブルーチンは図5-4-1(1)の②でコールされ、ICCG法の収束性を高めるための前処理として図5-4-5(1)のように  $A \doteq L D L^T$  の分解を行います。これを不完全  $L D L^T$  分解と呼びます。分解の程度によって ICCG(1,1)、ICCG(1,2)などの方法があり(参考文献[22]参照)、ICCG法の収束性に影響しますが、ここでは最も簡単なICCG(1,1)で分解します。分解された  $L D L^T$  を4本の帯  $A_i, B_i, C_i, D_i$  で表すことにします。このとき  $D_i = 1/A_i$  となるようにします。

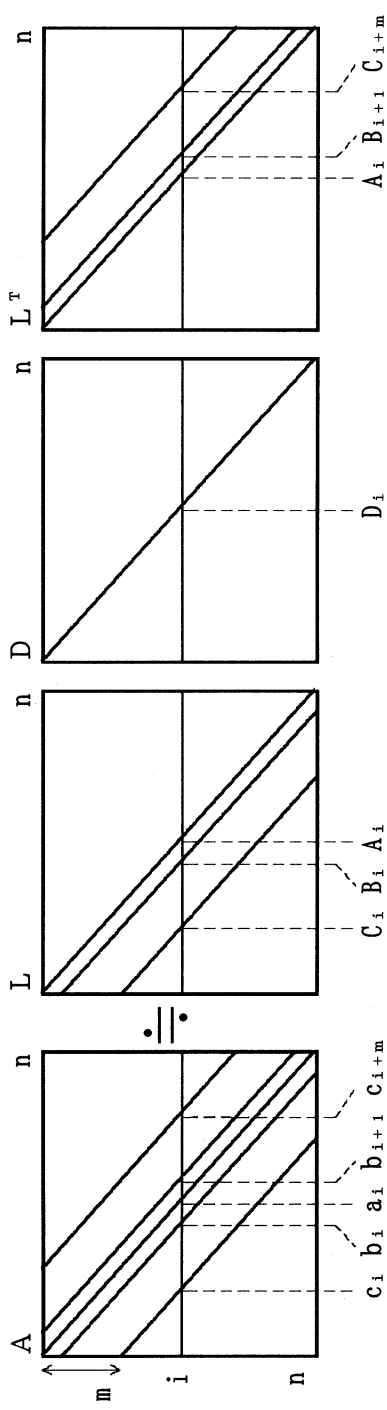


図5-4-5(1)

ICCG(1,1)の場合、行列  $A$  の  $i$  行目の5つの各要素と、 $L \cdot D \cdot L^T$  を実際に掛けた行列の同じ位置の要素の み が等しくなるようにします。ちなみにこれ以外の要素については両辺で異なってもよいので、この分解方法を不完全な  $L D L^T$  分解と呼びます。その結果以下の①の関数式が導かれます。この式から分かるように、分解後の  $B_i, C_i$  は分解前の  $b_i, c_i$  と一致するので実際には求める必要はありません。また①の3本目の式の  $A_i$  を左辺にし、 $D_i = 1/A_i, B_i = b_i, C_i = c_i$  を使用して整理すると図式となり、これを前述のように2次元配列で表すと図式となります。本サブルーチンでは結局③式から対角要素  $D_{i,j}$  だけを求めます。

③式で配列  $D$  に着目すると図5-4-5(2)のようになり、配列  $D$  には「左上の要素から自分を求める」という依存関係があるため完全な並列性はありませんが、このパターンは4-6-7節で説明したパイプライン法で並列化することができます。ただし図5-4-1(1)の②から分かるように、本サブルーチンは収束反復のループに入る前に1回しかコールされず通常ホットスポットにはならないので、今回は並列化を行わず、図5-4-10の③に示すように全プロセスが配列  $D$  の全ての要素を求めるようにしました(もし並列化する場合は後述するサブルーチン  $L D L^T$  と同じ処理になります)。なお、図5-4-1(2)、図5-4-10の③の番号と本文内の③の番号は対応しています。

余談ですが、図5-4-5(2)のうち着色した部分は左の要素からの依存関係がないので、本来は③式を場合分けする必要があります。ところが場合分けすると、サブルーチン  $L D L^T$  をパイプライン法で並列化する際に処理が面倒になるので、ここでは図5-4-5(3)のように周囲の1要素分にゼロを入れ、全要素に対して③式をそのまま適用するようになりました。

$$\begin{cases} C_i = C_i \\ b_i = B_i \\ a_i = D_{i-1} B_i^2 + D_{i-m} C_i^2 + A_i \\ b_{i+1} = B_{i+1} \\ C_{i+m} = C_{i+m} \end{cases} \quad \text{①}$$

$$\begin{cases} D_0 \quad i = 1, n \\ D_i = 1 / (a_i - D_{i-m} c_i^2 - D_{i-1} b_i^2) \dots \\ \text{ENDDO} \end{cases} \quad \text{②}$$

$$\begin{cases} D_0 \quad j = 1, M Y \\ D_0 \quad i = 1, M X \\ D_{i,j} = 1 / (a_{i,j} - D_{i,j-1} c_{i,j}^2 - D_{i-1,j} b_{i,j}^2) \dots \\ \text{ENDDO} \\ \text{ENDDO} \end{cases} \quad \text{③}$$

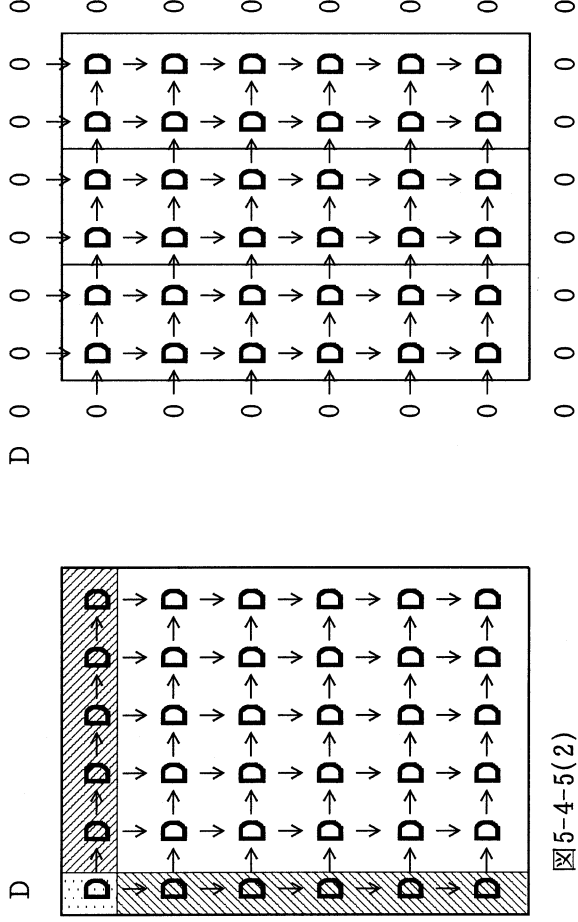


図5-4-5(2)

図5-4-5(3)

■ サブルーチンAXSUB

本サブルーチンは図5-4-1(1)の③,⑦でコールされ、行列ベクトル積  $y = Ax$  を計算します。ここでは解ベクトルではなく、 $x$ と $y$ は任意のベクトルを表します。図5-4-6(1)の*i*行目の関係式から④式が得られ、これを2次元配列で表すと⑤式となります。

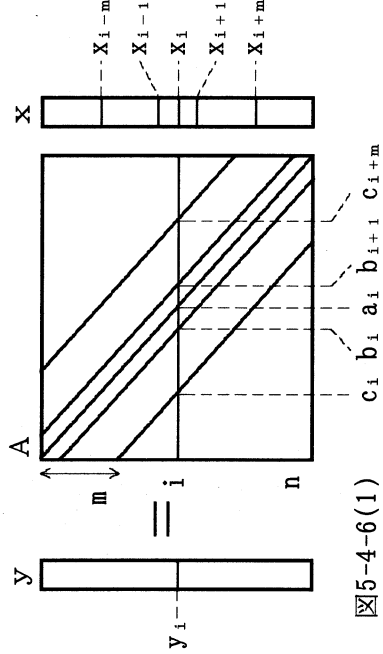


図5-4-6(1)

$$\begin{aligned}
 & \text{DO } i = 1, n \\
 & \quad y_i = c_i x_{i-m} + b_i x_{i-1} + a_i x_i + b_{i+1} x_{i+1} + c_{i+m} x_{i+m} \quad \text{④} \\
 & \text{ENDDO}
 \end{aligned}$$

```

DO j = 1, MY
DO i = 1, MX
  y_{i,j} = c_{i,j} x_{i,j-1} + b_{i,j} x_{i,j} + a_{i,j} x_{i,j} + b_{i+1,j} x_{i+1,j} + c_{i,j+m} x_{i,j+m} \quad \text{⑤}
ENDDO
ENDDO
  
```

⑤式で配列 $x$ と $y$ に着目すると図5-4-6(2)のようになり、配列 $y$ のある要素 $Y$ は配列 $x$ の同じ位置と上下左右の要素 $X$ を使用して計算されます。配列 $y$ の要素間には依存関係がない(つまりどの順に計算してもよい)ので⑤式のループは並列化することができます。ただし、他の部分の並列化によって各プロセスは配列 $x$ のうち自分の担当部分しか持っていないため、⑤式の計算に入る前に図5-4-6(3)に示す非循環シフトで境界の1列を隣のプロセスとの間で交換する必要があります(4-6-2節参照)。これを図5-4-10の④で行います。なお、⑤式の $a, b, c$ の各配列は全プロセスが全ての値を持っているので、計算前に通信を行う必要はありません。

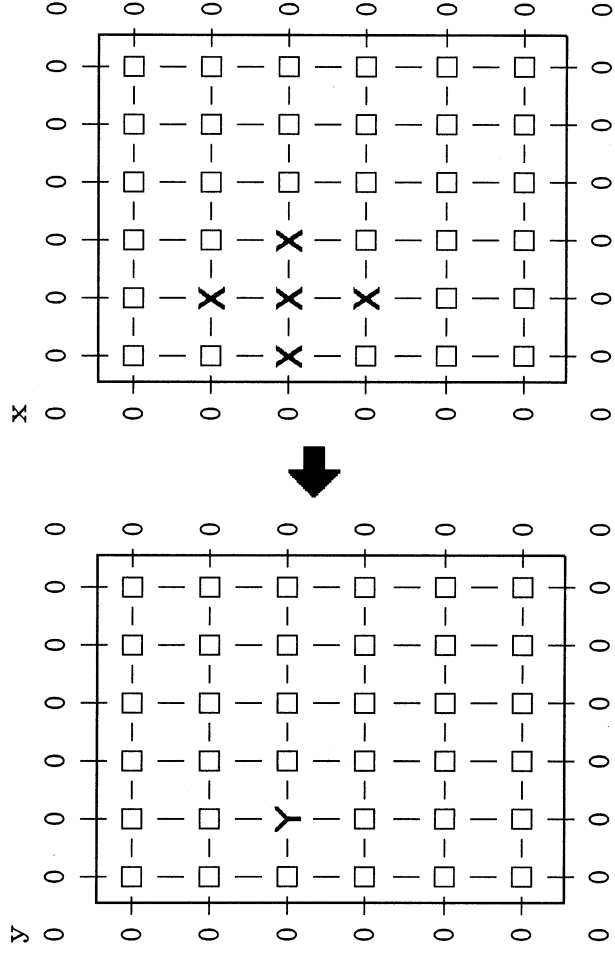


図5-4-6(2)

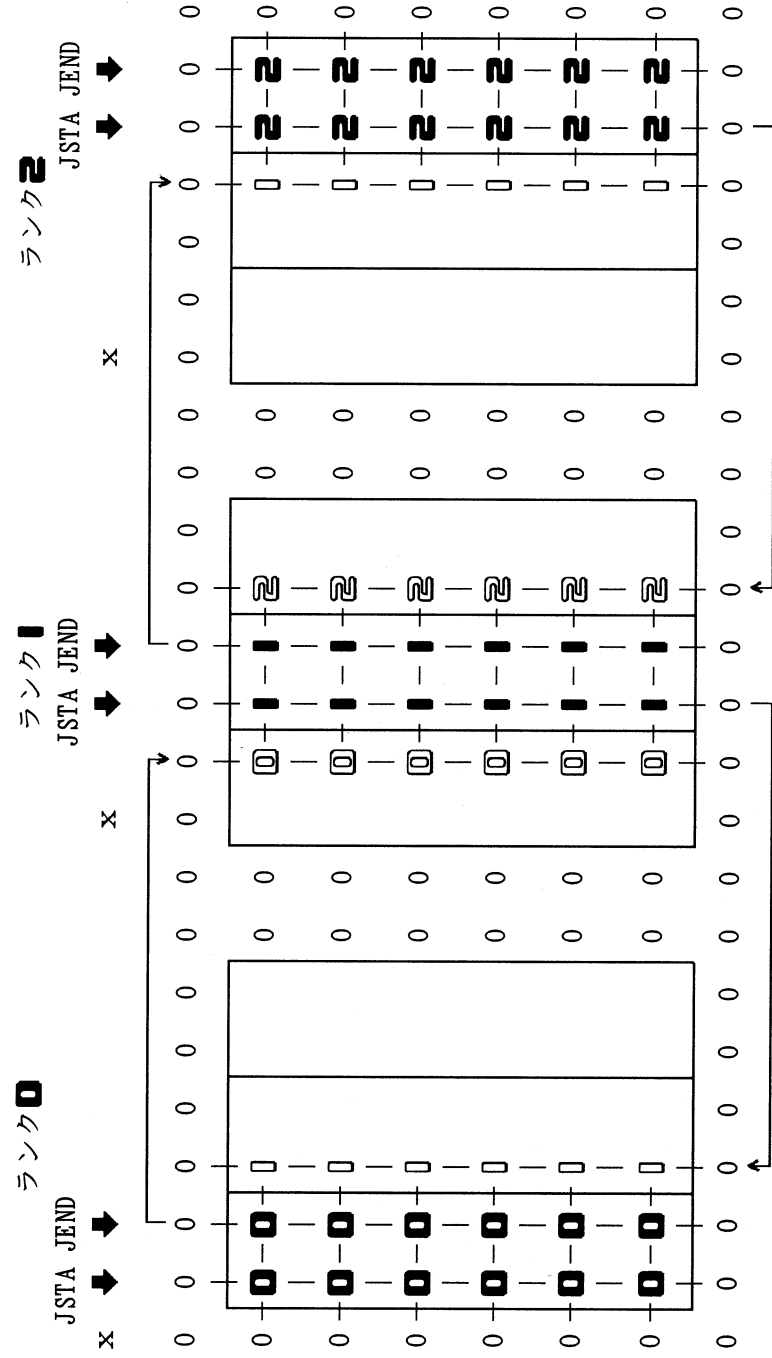


図5-4-6(3)



■ サブルーチン  $LDL^T$

本サブルーチンは図5-4-1(1)の⑤,⑬でコールされ、 $y = (LDL^T)^{-1}x$  の計算を行ないます。ここで  $x$  は解ベクトルではなく、 $x$  と  $y$  は任意のベクトルを表します。この演算を以下のように変形して作業ベクトル  $w$  を導入すると、前進消去と後退代入の2つのステップに分解されます。

$$y = (LDL^T)^{-1}x = (DL^T)^{-1}L^{-1}x$$

↳  $w$

【前進消去】  $w = L^{-1}x$  から  $w$  を求める。 ➡ 実際には  $x = Lw$  から  $w$  を求めます。

【後退代入】  $y = (DL^T)^{-1}w$  から  $y$  を求める。 ➡ 実際には  $w = DL^T y$  から  $y$  を求めます。

まず前進消去のステップでは、図5-4-7(1)に示すように  $i$  行目の関係式は

$$x_i = c_i w_{i-m} + B_i w_{i-1} + A_i w_i$$

となります。  $w_i$  を左辺にし、 $C_i = c_i$ ,  $B_i = b_i$ ,  $A_i = 1/D_i$  を使用して整理すると⑥式が得られます。

一方後退代入のステップでは、図5-4-7(2)に示すように  $i$  行目の関係式は

$$w_i = A_i D_i y_i + B_{i+1} D_i y_{i+1} + C_{i+m} D_i y_{i+m}$$

となります。  $y_i$  を左辺にし、 $C_i = c_i$ ,  $B_i = b_i$ ,  $A_i = 1/D_i$  を使用して整理すると⑦式が得られます。

実際のプログラムでは、⑥,⑦式の作業ベクトル  $w_i$  を、最終結果を入れるベクトル  $y_i$  で代用し、⑧,⑨式で計算します。これを2次元配列で表すと⑩,⑪式のようになります。

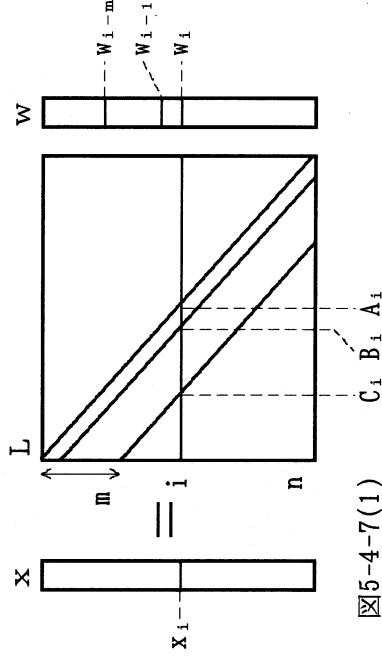


図5-4-7(1)

```

DO i = 1, n
  w_i = D_i * (x_i - c_i w_{i-m} - b_i w_{i-1}) ... ⑥
ENDDO
DO i = n, 1, -1
  y_i = w_i - D_i * (b_{i+1} y_{i+1} + c_{i+m} y_{i+m}) ... ⑦
ENDDO

```

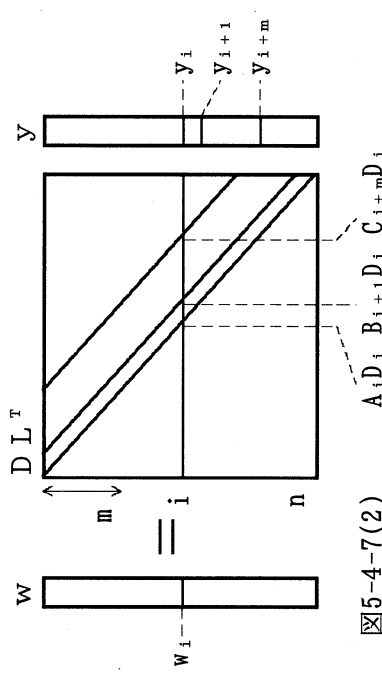


図5-4-7(2)

```

DO i = 1, n
  y_i = D_i * (x_i - c_i y_{i-m} - b_i y_{i-1}) ... ⑧
ENDDO
DO i = n, 1, -1
  y_i = y_i - D_i * (b_{i+1} y_{i+1} + c_{i+m} y_{i+m}) ... ⑨
ENDDO

```

```

DO j = 1, MY
  DO i = 1, MX
    y_{i,j} = D_{i,j} * (x_{i,j} - c_{i,j} y_{i,j-1} - b_{i,j} y_{i-1,j}) ... ⑩
  ENDDO
ENDDO
DO j = MY, 1, -1
  DO i = MX, 1, -1
    y_{i,j} = y_{i,j} - D_{i,j} * (b_{i+1,j} y_{i+1,j} + c_{i,j+1} y_{i,j+1}) ... ⑪
  ENDDO
ENDDO

```

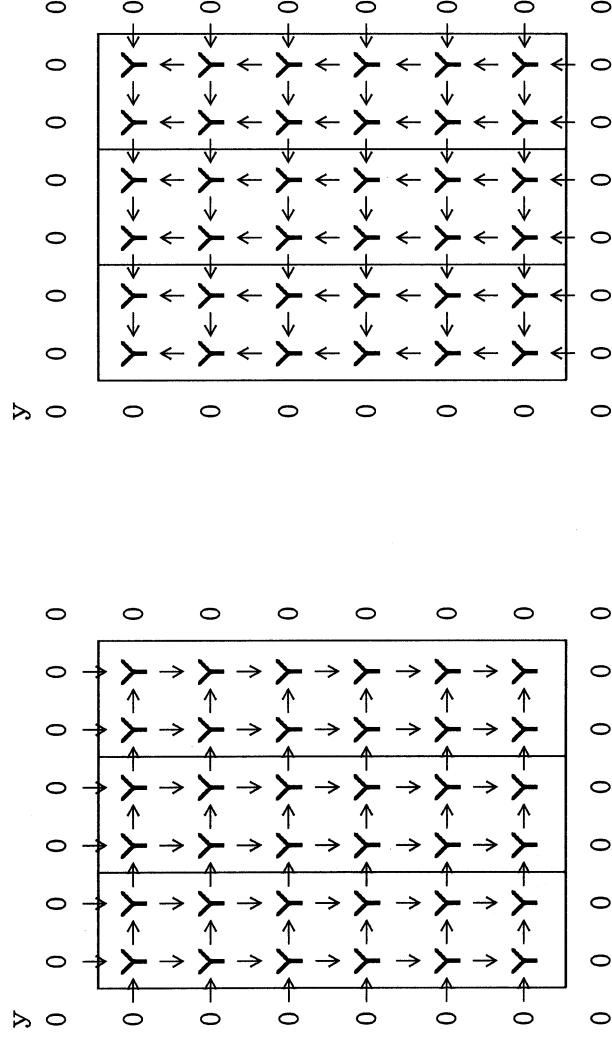


図5-4-8(1)

前進消去を行う図式で配列yに着目すると図5-4-8(1)のようになり、配列yには「左と上の要素から自分を求める」という依存関係があるため完全な並列性はありませんが、このパターンは4-6-7節で説明したパイプライン法で並列化することができます。

まず配列yを1次元方向に大きさIBLOCK(図5-4-8(3)の $\updownarrow$ )のブロックに分割します。IBLOCKの大きさは図5-4-10の[2]で指定します(最適な大きさは試行錯誤で調整して下さい)。各プロセスの計算順序は図5-4-8(3)のようになり、□内の数字の順に処理を行い、○内の数字の順に通信を行います。これを図5-4-10の[4]と[4]で行います(詳細は4-6-7節を参照して下さい)。 $\square$ のIBLKLENは処理するブロックの1次元方向の大きさを基本的にはIBLOCKですが、MXがIBLOCKで割りきれない場合、最後の半端な部分ではIBLOCKより小さくなるのでそれを調整します。なお、 $\square$ 式のb,c,dの各配列は全プロセスが全ての値を持っており、配列Xは各プロセスが自分の担当部分の値を持っているので計算の前には通信する必要はありません。

一方後退代入を行う図式で配列yに着目すると図5-4-8(2)のようになり、配列yに図5-4-8(1)と逆の依存関係があるので、図5-4-8(4)のような順序で処理を行います。これを図5-4-10の[4]と[7]、[8]で行います。 $\square$ のIBLKLENは前述の[4]と同じで、計算領域の1次元方向が『1』から開始している場合はこのようになりますが、例えば『2』から開始している場合は  $IBLKLEN = \text{MIN}(IBLOCK, \text{II}-1)$  となるので注意して下さい。

ランク0

JSTA JEND

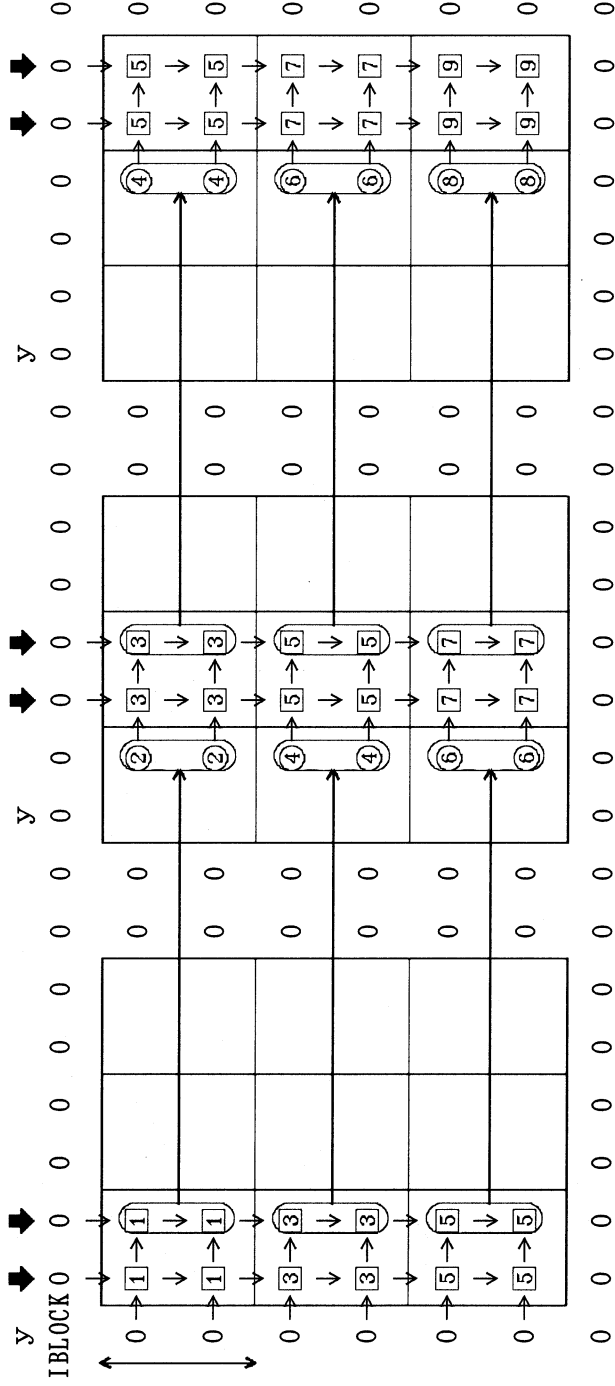
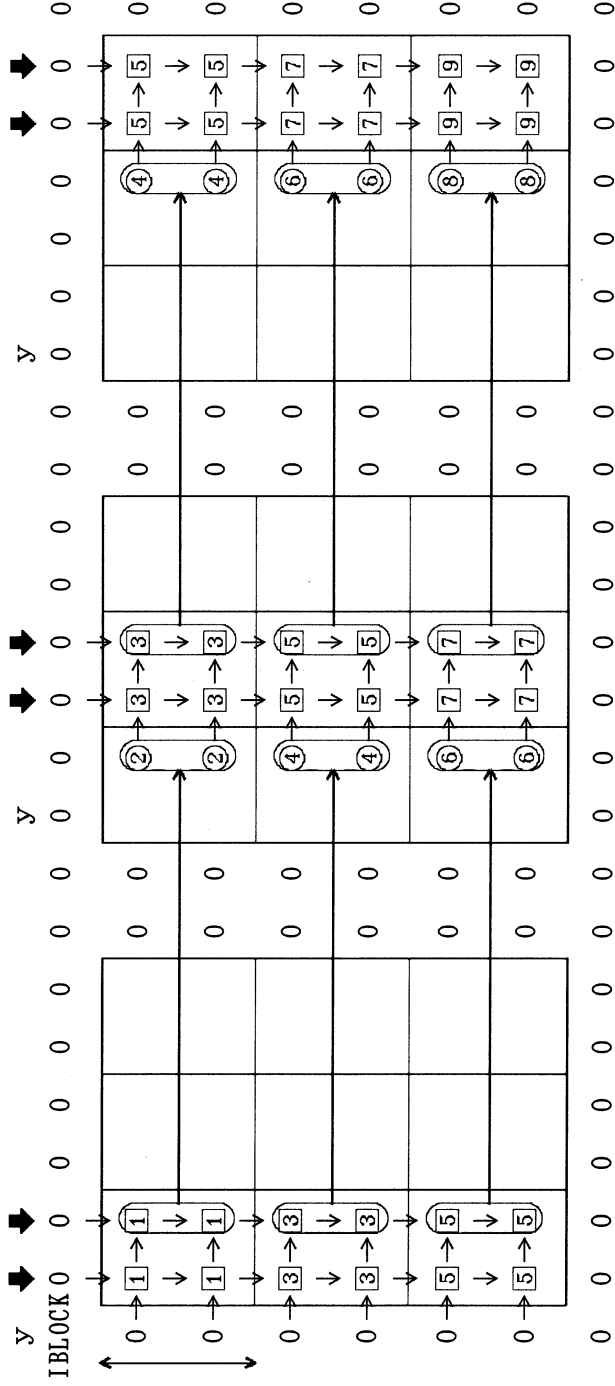


図5-4-8(3)

ランク1

JSTA JEND



ランク0

JSTA JEND

```

y  ↓ ↓
0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
    ↓ ↓

```

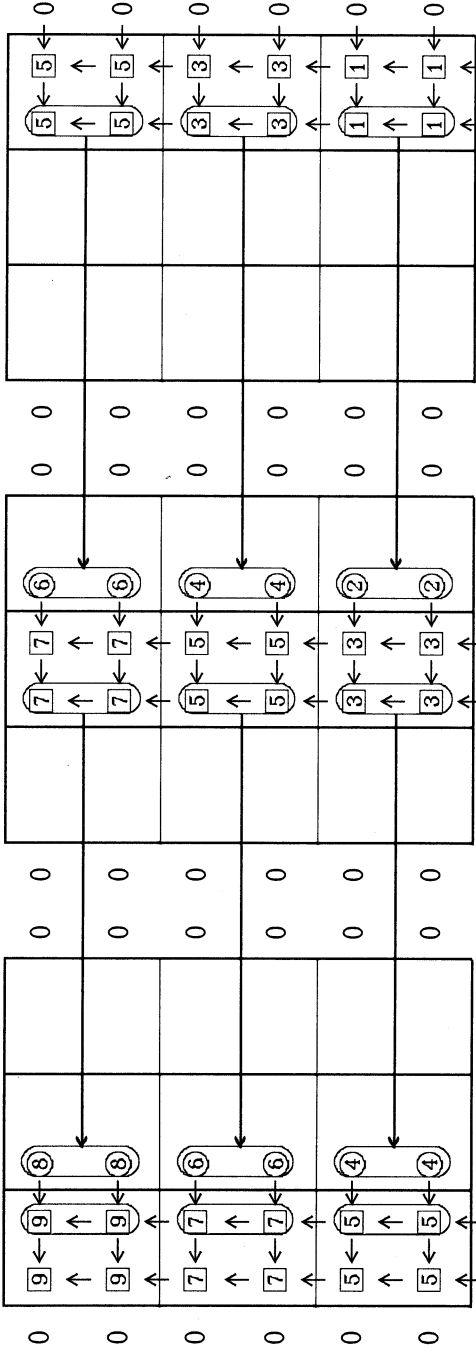


図5-4-8(4)

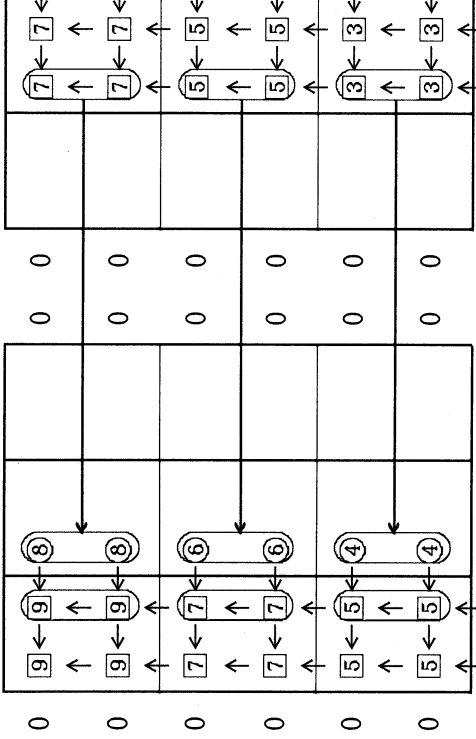
ランク1

JSTA JEND

```

y  ↓ ↓
0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
    ↓ ↓

```



① ↷

```

MODULE PARA
  INCLUDE 'mpif.h'
  PARAMETER (NCPUs=3)
  INTEGER NPROCS,MYRANK,JSTA,JEND
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  INTEGER IUP, IDOWN, ITYPE(0:NCPUs-1)
  END

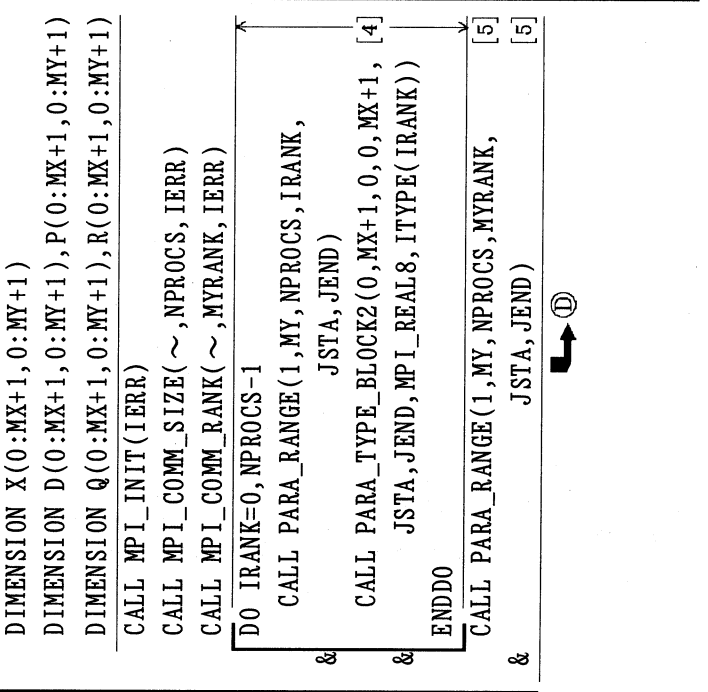
PROGRAM MAIN
  USE PARA
  IMPLICIT REAL*8(A-H,0-Z)
  PARAMETER (MX=6,MY=6)
  DIMENSION A(0:MX+1,0:MY+1,3),B(MX,MY)
  DIMENSION X(0:MX+1,0:MY+1)
  DIMENSION D(0:MX+1,0:MY+1),P(0:MX+1,0:MY+1)
  DIMENSION Q(0:MX+1,0:MY+1),R(0:MX+1,0:MY+1)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~,NPROCS,IERR)
  CALL MPI_COMM_RANK(~,MYRANK,IERR)
  DO IRANK=0,NPROCS-1
    CALL PARA_RANGE(1,MY,NPROCS,IRANK,
      JSTA,JEND)
  ENDDO
  &
  &
  &
  CALL PARA_RANGE(1,MY,NPROCS,MYRANK,
    JSTA,JEND)
  ENDDO

```

```

IUP = MYRANK+1
IDOWN = MYRANK-1
IF (IUP==NPROCS) IUP = MPI_PROC_NULL
IF (IDOWN==-1) IDOWN = MPI_PROC_NULL
:
配列AとBに値を設定します。
EPS=1.0D-7
ITR=1000
CALL ICCG(A,MX,MY,B,EPS,ITR,X,D,P,Q,R,IER)[3]
:
CALL MPI_FINALIZE(IERR)
END
SUBROUTINE ICCG(A,MX,MY,B,EPS,ITR,X,D,P,Q,R,IER)
& X,D,P,Q,R,IER)
USE PARA
IMPLICIT REAL*8(A-H,0-Z)
DIMENSION A(0:MX+1,0:MY+1,3),B(MX,MY)
DIMENSION X(0:MX+1,0:MY+1)
DIMENSION D(0:MX+1,0:MY+1)
DIMENSION Q(0:MX+1,0:MY+1),R(0:MX+1,0:MY+1)
DIMENSION WORKS(2),WORKR(2)
DO J = 0, MY+1
  DO I = 0, MX+1
    D(I,J) = 0.0D0
    X(I,J) = 0.0D0
    P(I,J) = 0.0D0
    Q(I,J) = 0.0D0
    R(I,J) = 0.0D0
  ENDDO
ENDDO

```



① ↷

② ↷

← ⊕

```

CALL DECOMP(A,D,MX,MY)
CALL AXSUB(A,X,Q,MX,MY)
DO J = JSTA,JEND
  DO I = 1, MX
    R(I,J) = B(I,J) - Q(I,J)
  ENDDO
ENDDO
CALL LDLT(A,D,R,P,MX,MY)
CC1 = 0.0D0
DO J = JSTA,JEND
  DO I = 1, MX
    CC1 = CC1 + R(I,J)*P(I,J)
  ENDDO
ENDDO
CALL MPI_ALLREDUCE(CC1,C1,1,MPI_REAL8,
  [6] MPI_SUM,MPI_COMM_WORLD,IERR) [6]
&
DO K = 1, ITR 収束反復ループ
CALL AXSUB(A,P,Q,MX,MY)
CC2 = 0.0D0
DO J = JSTA,JEND
  DO I = 1, MX
    CC2 = CC2 + P(I,J)*Q(I,J)
  ENDDO
ENDDO
CALL MPI_ALLREDUCE(CC2,C2,1,MPI_REAL8, [7]
  [7] MPI_SUM,MPI_COMM_WORLD,IERR) [7]
&
ALPHA = C1 / C2
X1 = 0.0D0
X2 = 0.0D0
DO J = JSTA,JEND
  DO I = 1, MX
    Y = X(I,J)
    X(I,J) = X(I,J) + ALPHA*P(I,J)
    R(I,J) = R(I,J) - ALPHA*Q(I,J)
    X1 = X1 + Y*Y
    X2 = X2 + (X(I,J)-Y)**2
  ENDDO
ENDDO

```

⊕ →

⊕ →

```

WORKS(1) = X1
WORKS(2) = X2
CALL MPI_ALLREDUCE(WORKS,WORKR,2,
  [8] MPI_REAL8,MPI_SUM,MPI_COMM_WORLD,IERR)
&
X1 = WORKR(1)
X2 = WORKR(2)
IF (X1 /= 0.0) THEN
  RES = DSQRT(X2/X1)
  IF (RES <= EPS) THEN
    ITR = K
    IER = 0
    GOTO 999
  ENDDO
  ENDDO
  CALL LDLT(A,D,R,Q,MX,MY)
  CC3 = 0.0D0
  DO J = JSTA,JEND
    DO I = 1, MX
      CC3 = CC3 + R(I,J)*Q(I,J)
    ENDDO
  ENDDO
  CALL MPI_ALLREDUCE(CC3,C3,1,MPI_REAL8, [9]
    [9] MPI_SUM,MPI_COMM_WORLD,IERR) [9]
  &
  BETA = C3/C1
  C1 = C3
  DO J = JSTA,JEND
    DO I = 1, MX
      P(I,J) = Q(I,J) + BETA*P(I,J)
    ENDDO
  ENDDO
  ENDDO
  IER = 1
  999 CONTINUE
  DO IRANK=0,NPROCS-1
    CALL MPI_BCAST(X,1,ITYPE(IRANK),IRANK,
      [10] MPI_COMM_WORLD,IERR) [10]
  ENDDO
  EPS = RES
  END

```

⊕ →



```
SUBROUTINE DECOMP(A,D,MX,MY)
```

```
USE PARA
```

```
IMPLICIT REAL*8(A-H,0-Z)
```

```
DIMENSION A(0:MX+1,0:MY+1,3)
```

```
DIMENSION D(0:MX+1,0:MY+1)
```

```
DO J = 1, MY
```

```
  DO I = 1, MX
```

```
    D(I,J) = 1.0D0/( A(I,J,3)
```

```
      -D(I,J-1)*A(I,J,1)**2
```

```
      -D(I-1,J)*A(I,J,2)**2 )
```

```
    ENDDO
```

```
  ENDDO
```

```
END
```

```
SUBROUTINE AXSUB(A,X,Y,MX,MY)
```

```
USE PARA
```

```
IMPLICIT REAL*8(A-H,0-Z)
```

```
DIMENSION A(0:MX+1,0:MY+1,3)
```

```
DIMENSION X(0:MX+1,0:MY+1),Y(0:MX+1,0:MY+1)
```

```
CALL MPI_ISEND(X(1,JEND),MX,MPI_REAL8,
```

```
& IUP , 1, MPI_COMM_WORLD, ISEND1, IERR)
```

```
CALL MPI_ISEND(X(1,JSTA),MX,MPI_REAL8,
```

```
& IDOWN, 1, MPI_COMM_WORLD, ISEND2, IERR)
```

```
CALL MPI_IRECV(X(1,JSTA-1),MX,MPI_REAL8,
```

```
& IDOWN, 1, MPI_COMM_WORLD, IRECV1, IERR) [14]
```

```
CALL MPI_IRECV(X(1,JEND+1),MX,MPI_REAL8,
```

```
& IUP , 1, MPI_COMM_WORLD, IRECV2, IERR)
```

```
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
```

```
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
```

```
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
```

```
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
```

```
DO J = JSTA, JEND
```

```
  DO I = 1, MX
```

```
    Y(I,J) = A(I,J,1)*X(I,J-1)
```

```
      + A(I,J,2)*X(I-1,J)
```

```
      + A(I,J,3)*X(I,J)
```

```
      + A(I+1,J,2)*X(I+1,J)
```

```
      + A(I,J+1,1)*X(I,J+1)
```

```
  ENDDO
```

```
ENDDO
```

```
END
```



```
SUBROUTINE LDLT(A,D,X,Y,MX,MY)
```

```
USE PARA
```

```
IMPLICIT REAL*8(A-H,0-Z)
```

```
DIMENSION A(0:MX+1,0:MY+1,3)
```

```
DIMENSION D(0:MX+1,0:MY+1)
```

```
DIMENSION X(0:MX+1,0:MY+1),Y(0:MX+1,0:MY+1)
```

```
IBLOCK=2
```

```
DO II=1,MX,IBLOCK
```

```
  IBLKEN = MIN( IBLOCK, MX-II+1) [13]
```

```
  CALL MPI_IRECV(Y(II,JSTA-1),IBLKLEN,
```

```
& MPI_REAL8, IDOWN, 1,
```

```
& MPI_COMM_WORLD, IREQR, IERR) [14]
```

```
  CALL MPI_WAIT(IREQR, ISTATUS, IERR)
```

```
  DO J=JSTA, JEND
```

```
    DO I=II, II+IBLKLEN-1
```

```
      Y(I,J) = D(I,J)*( X(I,J)
```

```
        -A(I,J,1)*Y(I,J-1)
```

```
        -A(I,J,2)*Y(I-1,J) )
```

```
    ENDDO
```

```
  ENDDO
```

```
  CALL MPI_ISEND(Y(II,JEND),IBLKLEN,
```

```
& MPI_REAL8, IUP, 1,
```

```
& MPI_COMM_WORLD, IREQS, IERR) [15]
```

```
  CALL MPI_WAIT(IREQS, ISTATUS, IERR)
```

```
ENDDO
```

```
DO II=MX, 1, -IBLOCK
```

```
  IBLKEN = MIN( IBLOCK, II) [16]
```

```
  III = II-IBLKLEN+1
```

```
  CALL MPI_IRECV(Y(III,JEND+1),IBLKLEN,
```

```
& MPI_REAL8, IUP, 1,
```

```
& MPI_COMM_WORLD, IREQR, IERR) [17]
```

```
  CALL MPI_WAIT(IREQR, ISTATUS, IERR)
```

```
  DO J=JEND, JSTA, -1
```

```
    DO I=II, II-IBLKLEN+1, -1
```

```
      Y(I,J) = Y(I,J)-D(I,J)*(
```

```
        A(I+1,J,2)*Y(I+1,J)
```

```
        +A(I,J+1,1)*Y(I,J+1) )
```

```
    ENDDO
```

```
  ENDDO
```

```
  CALL MPI_ISEND(Y(III,JSTA),IBLKLEN,
```

```
& MPI_REAL8, IDOWN, 1,
```

```
& MPI_COMM_WORLD, IREQS, IERR) [18]
```

```
  CALL MPI_WAIT(IREQS, ISTATUS, IERR)
```

```
ENDDO
```

```
END
```

図5-4-10

## 5-4-2 パラレル・ブロック・オーダリング法による並列化

本節では対称帯行列のICCG法を並列化する別の方法として、参考文献[42][39]で紹介されているパラレル・ブロック・オーダリング法(以後PB0法)について説明します。

### ■ 配列の分割方法

PB0法でもバイブライン法の場合と同様に、図5-4-1(1)の各ベクトルの2次元配列(計算領域と同じ形状)として取り扱います。以後、計算領域として図5-4-11(1)の $12 \times 12$ の2次元配列を想定します(バイブライン法の場合と同じく、配列の周囲が1要素分大きくなっています)。PB0法では並列化のための分割方法として、計算領域を1次元目と2次元目の両方でブロック分割します(5-1-3節参照)。図5-4-11(1)の例では1次元目を3個、2次元目を3個の計9プロセスで分割し、各プロセスは $4 \times 4$ の領域(以後ブロック)を担当します。

PB0法の並列版のプログラムの図5-4-16に示します(図5-4-1(1)内の番号と図5-4-16内の番号は対応しています)。まず[1]で1次元目と2次元目のプロセス数であるIPROCSとJPROCSを指定し、次に各プロセスの上下および左右のプロセスのランクを知るため、[2]と[3]で図5-4-11(2)のようなテーブルITABLEを作成し、[4]で自分の1次元と2次元方向のランクをMYRANKI, MYRANKJに設定します。ランク4のプロセスのMYRANKIとMYRANKJを図5-4-11(2)に示します。次に[5]で、各プロセスが担当する1次元方向と2次元方向の担当範囲をISTA, IEND, JSTA, JENDに求めます(ランク4のブロックの担当範囲を図5-4-11(1)に示します)。なお、説明を簡単にするため、配列の縮小は行なっていません。

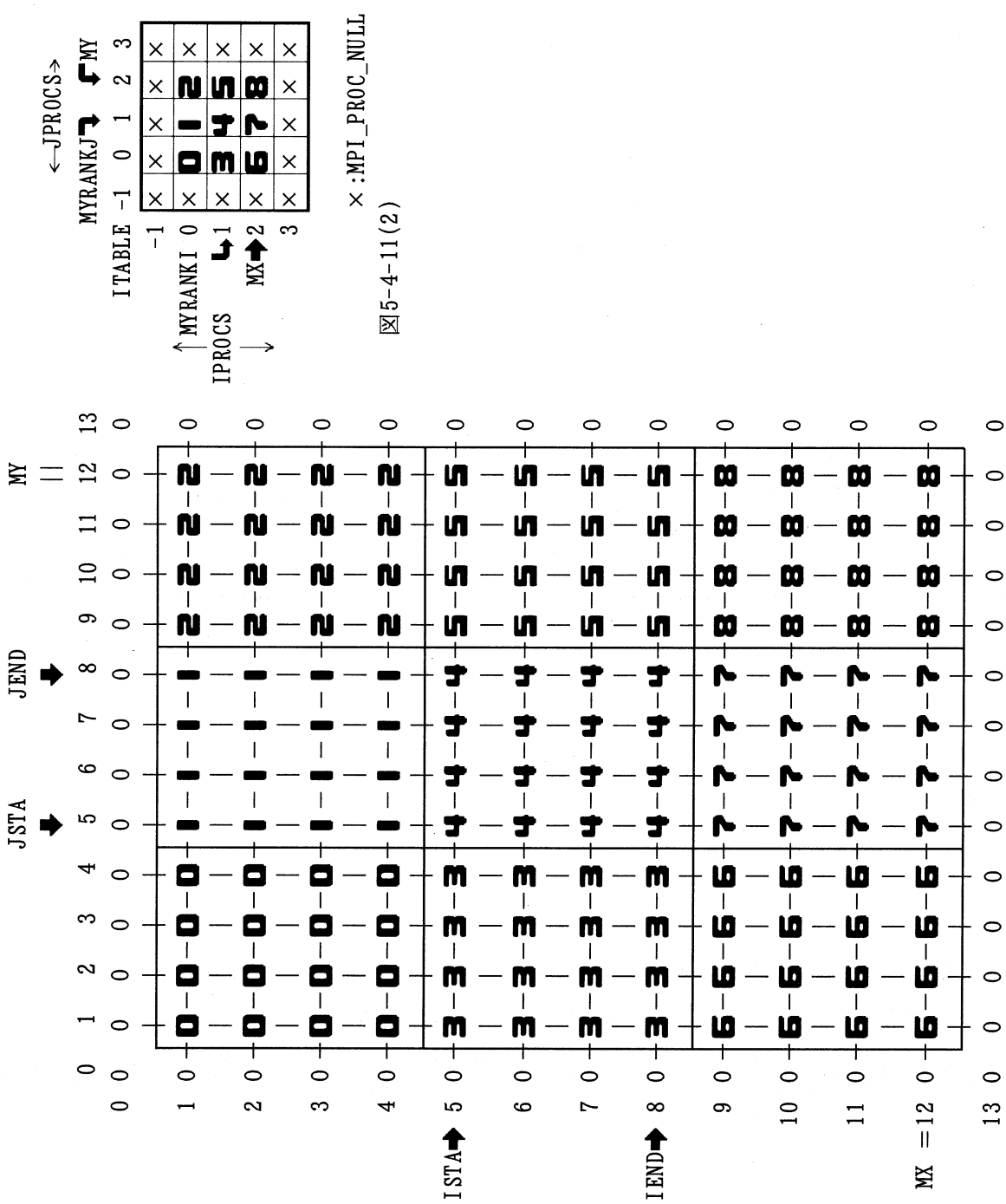


図5-4-11(1)

**ベクトル同志の演算と内積計算の並列化**

図5-4-1(1)の④,⑥,⑧,⑩,⑫,⑭,⑯の並列化にもなって必要になる[6]~[9]の通信と、最後に各プロセスが求めた解ベクトル $x$ の値を全プロセスに送信する[10],[11]の通信は、パイプライン法の場合と同じなので説明は省略します。

なお、[11]では3-5-3-2節で説明した自作のサブルーチン**PARA\_TYPE\_BLOCK2**を使用しましたが、MPI-2が使用できるマシン環境の場合は、MPI-2で提供されている**MPI\_TYPE\_CREATE\_SUBARRAY**(3-5-3-1節参照)を使用して下さい。

**サブルーチン D E C O M P**

PB0法では、例えば図5-4-12(1)の計算領域を図5-4-12(2)のようにブロック分割した場合、要素を図5-4-12(2)内の数字の順に並べて図5-4-11(3)の連立一次方程式を作成します(通常のICCG法の図5-4-2(4)に相当)。この連立一次方程式から前節の③,⑤,⑩,⑪式に相当する関係式を求めますが、導出過程は前節と同様なので省略し、以後の説明では導出した結果をいきなり使用します。

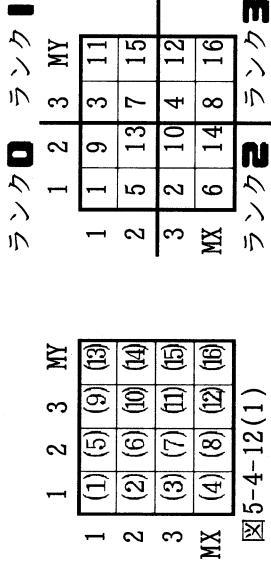
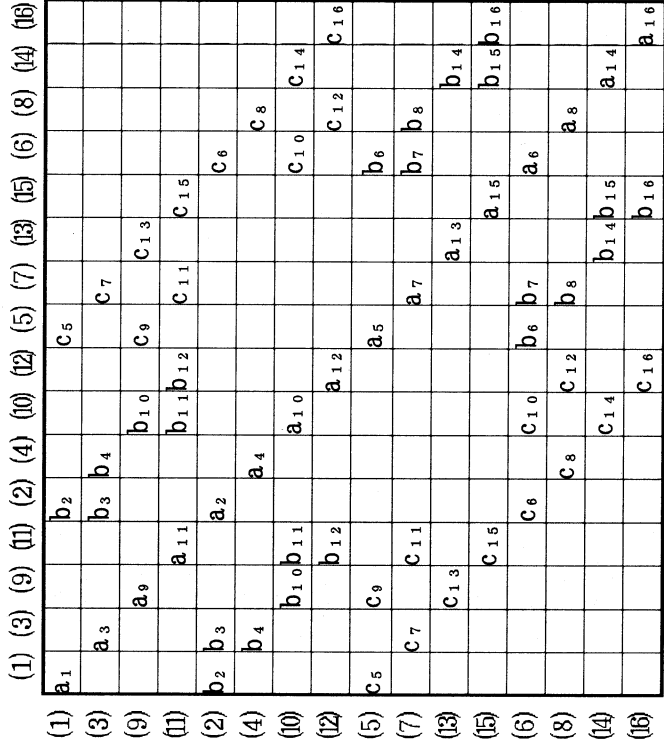


図5-4-12(2)



本サブルーチンでは  $A \equiv L D L^T$  を行いますが、通常のICCG法と同様に実際には対角要素  $D_{i,j}$  のみを求めます。通常のICCG法では図5-4-5(2)のように「左と上の要素から自分を求める」という依存関係があります。PB0法でも図5-4-13(1)のように、基本的には「左と上の要素から自分を求める」という依存関係がありますが、各ブロックの下の1行(7, 8, 9)では下の要素からの依存関係(↑)が、右の1列(9, 10, 13)では右の要素からの依存関係(←)があります。図中の□は同じ依存関係の要素を示します。このような依存関係があるので、□の数字の小さい要素から順に計算を行います(勿論、図からわかるように例えば□と□はどちらを先に計算してもかまいません)。

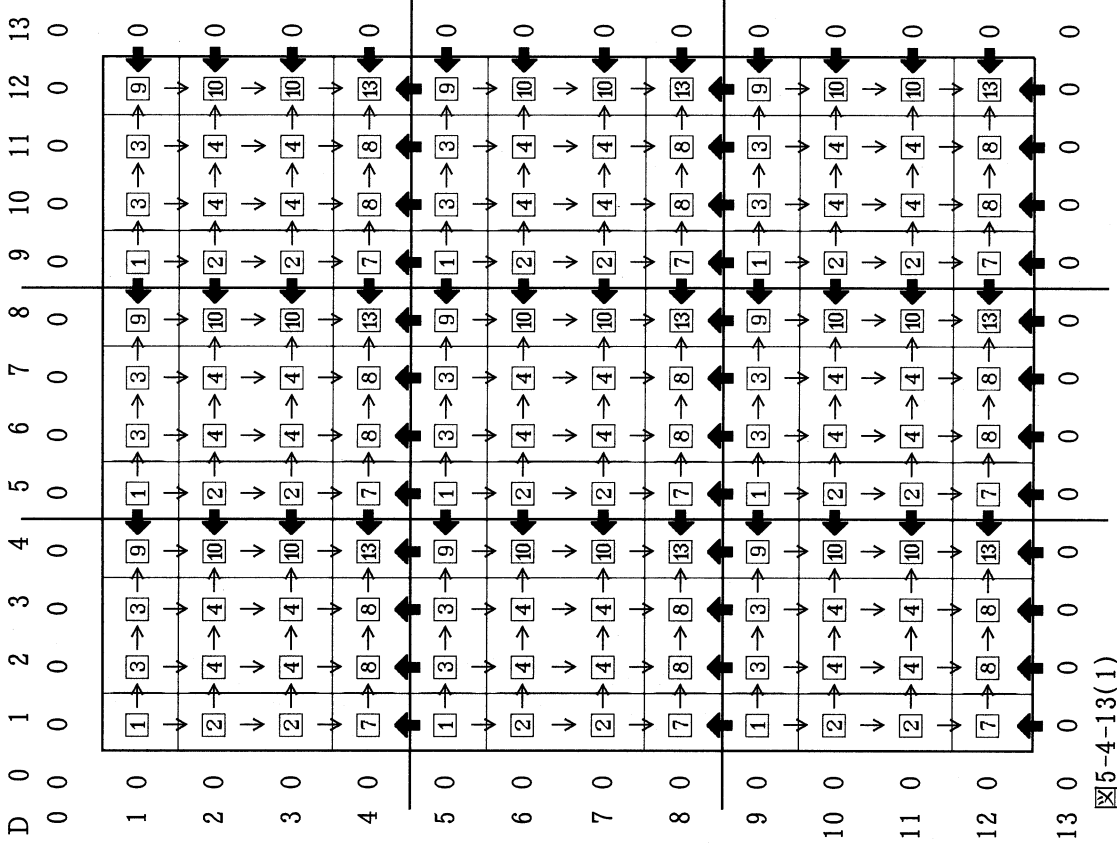


図5-4-13(1)

各要素の具体的な計算方法をまとめると図5-4-13(2)のようになります。基本的には前節の図式と同じですが、依存関係によって項の数が変わります。例えば図の要素の依存関係を(→↓↑←)に示しますが、このとき下線部の→, ↓, ↑, ←, ◀の項が必要になります。

$$\begin{aligned}
 \text{①} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j}) \\
 \text{②} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i-1,j} b_{i,j}^2) \\
 \text{③} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i,j-1} c_{i,j}^2) \\
 \text{④} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i,j-1} c_{i,j}^2 - D_{i-1,j} b_{i,j}^2) \\
 \text{⑦} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i-1,j} b_{i,j}^2 - D_{i+1,j} b_{i+1,j}^2) \\
 \text{⑧} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i,j-1} c_{i,j}^2 - D_{i-1,j} b_{i,j}^2) \\
 \text{⑨} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i,j-1} c_{i,j}^2 - D_{i-1,j} b_{i,j}^2 - D_{i,j+1} c_{i,j+1}^2) \\
 \text{⑩} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i,j-1} c_{i,j}^2 - D_{i-1,j} b_{i,j}^2 - D_{i,j+1} c_{i,j+1}^2) \\
 \text{⑬} & \left( \begin{array}{c} \phantom{\downarrow} \\ \phantom{\rightarrow} \\ \phantom{\leftarrow} \\ \phantom{\uparrow} \end{array} \right) D_{i,j} = 1/(a_{i,j} - D_{i,j-1} c_{i,j}^2 - D_{i-1,j} b_{i,j}^2 - D_{i+1,j} b_{i+1,j}^2 - D_{i,j+1} c_{i,j+1}^2)
 \end{aligned}$$

図5-4-13(2)



次に並列化した場合の動作について図5-4-13(3)で説明します。図5-4-16では配列の縮小を行っていませんが、図5-4-13(3)では紙面の関係で各プロセスが担当する部分の周辺のみを示します(図の周囲については座標番号に注意して下さい)。例えば図5-4-13(3)でランク4のプロセスに着目すると、処理は以下のようなになります。実際のプログラムは図5-4-16のサブルーチンDECOMPの対応する□を参照して下さい。

- ①～④の要素をこの順に処理します。
- ⑤と⑥の要素を左上のプロセスに対して送信し、右と下のプロセスから受信します。これによって⑦～⑩の要素の処理が可能になります。
- ⑦～⑩の要素をこの順に処理します。
- ⑪と⑫の要素を左上のプロセスに対して送信し、右と下のプロセスから受信します。これによって⑬の要素の処理が可能になります。
- ⑬の要素を処理します。

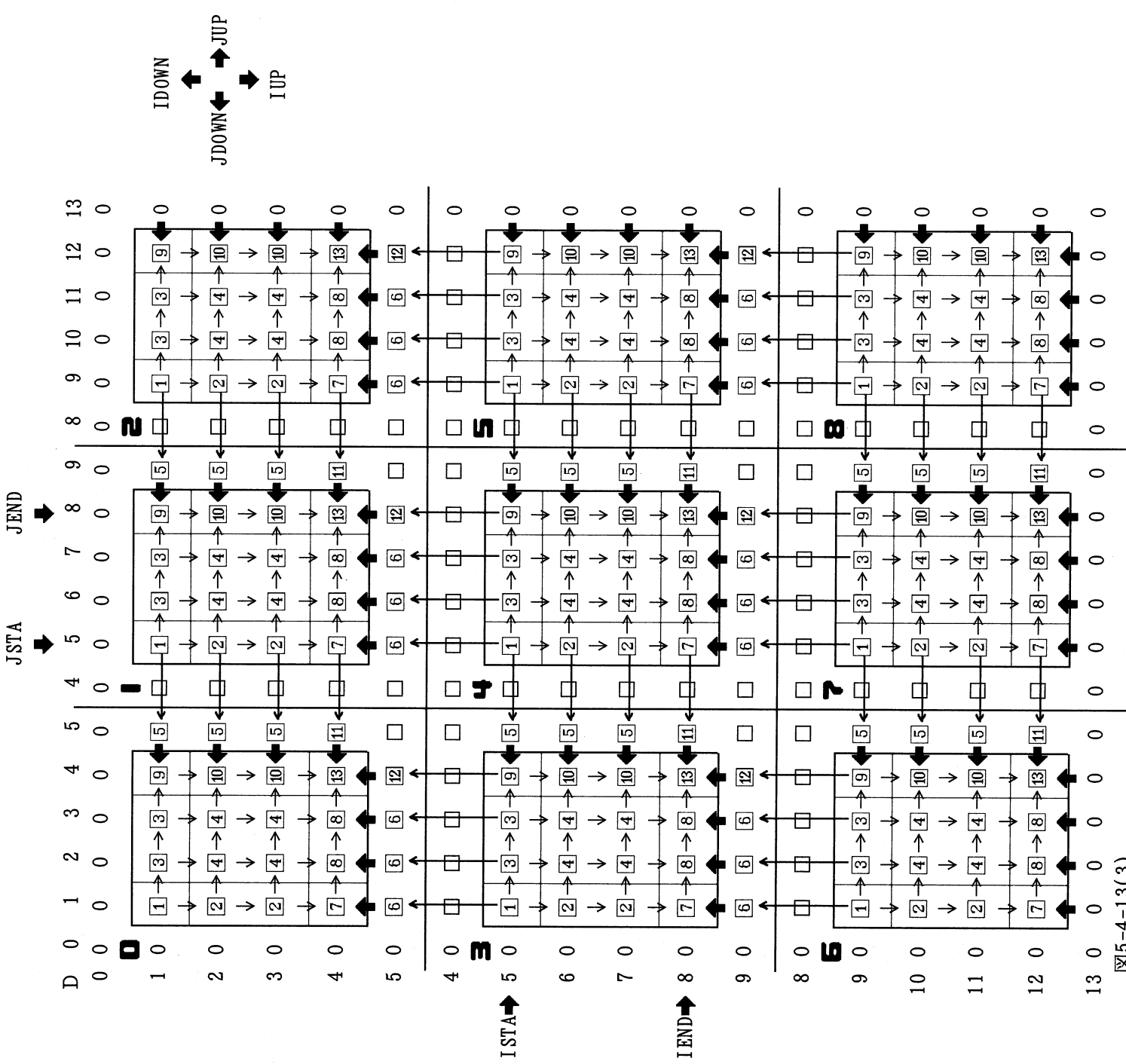


図5-4-13(3)

■ サブルーチン **AXSUB**

本サブルーチンの計算方法は前節の図式と同じで以下ようになります。PB0法では計算領域を1次元目と2次元目の両方で分割しているため、図5-4-6(3)に対応する境界のデータの通信は、図5-4-14の矢印のように上下左右のプロセスに対して行います。計算部分と通信部分を図5-4-16の図、図に示します。

```

D0 j = 1, MY
  D0 i = 1, MX
    Yi,j = Ci,jXi,j-1 + bi,jXi,j + ai,jXi,j+1 + bi+1,jXi+1,j + Ci,j+1Xi,j+1 ...
  ENDDO
ENDDO

```

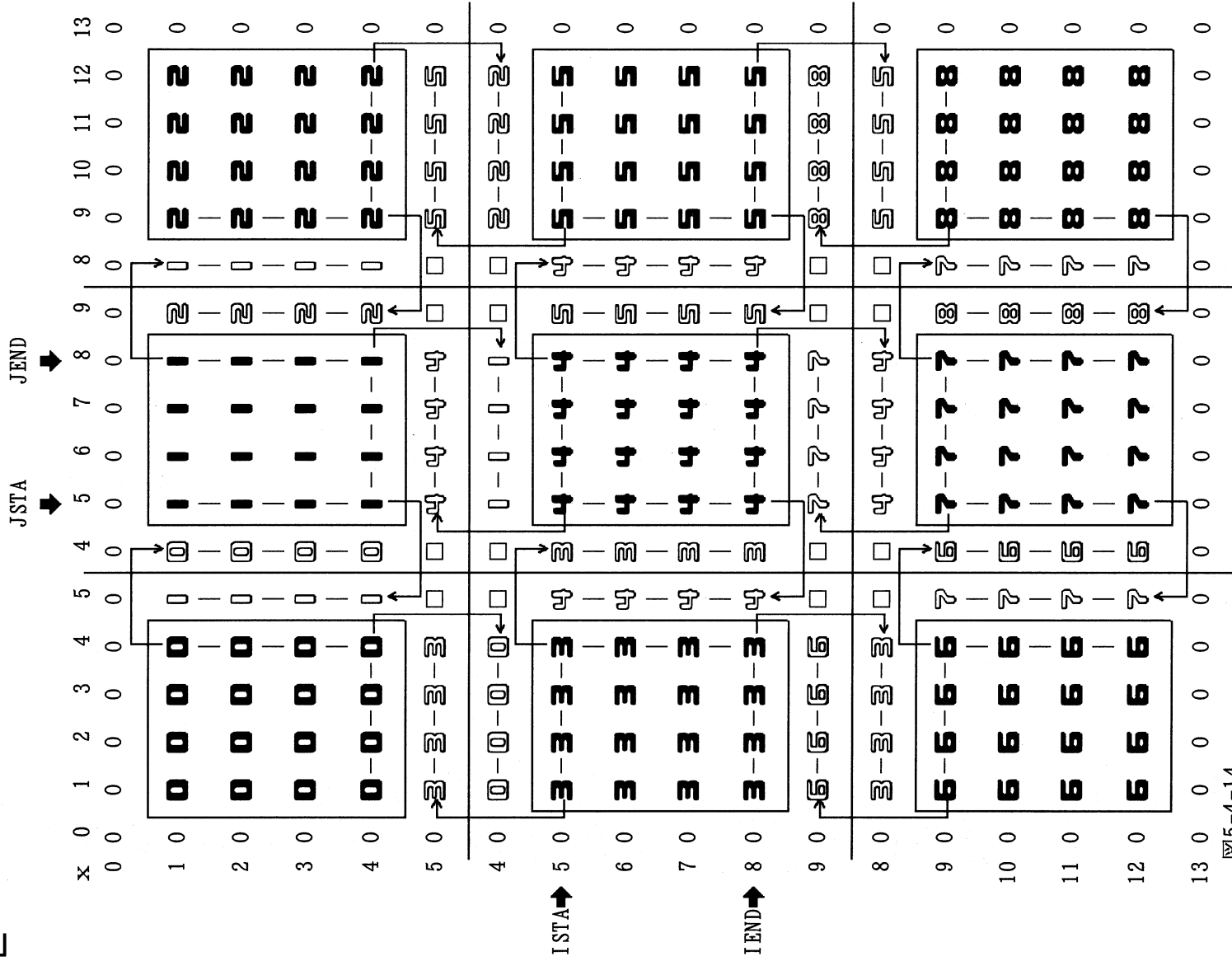


図5-4-14

■ サブルーチンLDLT

前節で説明したように、本サブルーチンには前進消去と後退代入の2つのステップがあり、このうち前進消去のステップでの並列化の方法はサブルーチンDECOMPと全く同じです。図5-4-8(1)に対応する要素間の依存関係は前述の図5-4-13(1)のようになり、それを並列化した場合の動作は図5-4-13(3)のようになります。各要素の具体的な計算方法は図5-4-15(1)のようになります。実際のプログラムは図5-4-16のサブルーチンLDLTの対応する□を参照して下さい。

【前進消去】

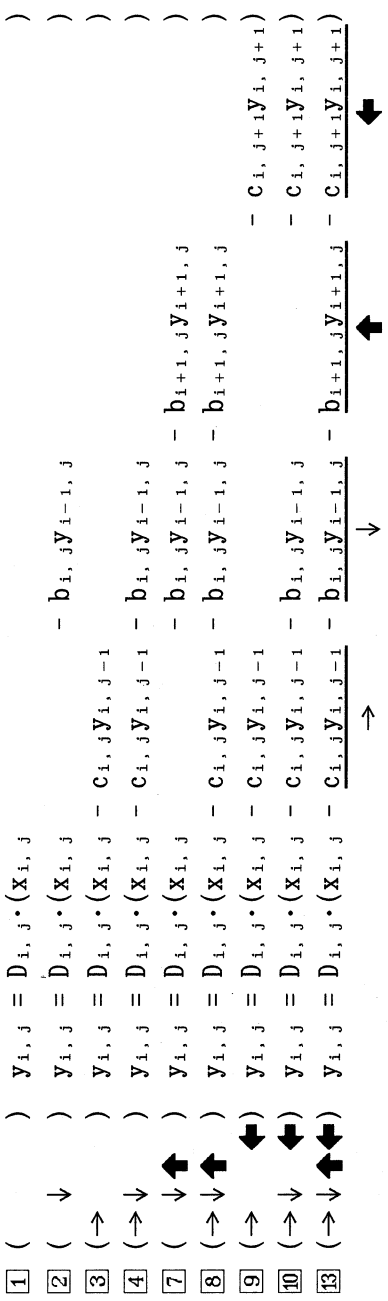


図5-4-15(1)

後退代入の動作は前進消去と逆で、並列化した場合は図5-4-15(3)のようになります。各要素の具体的な計算方法をまとめると図5-4-15(2)のようになります。実際のプログラムは図5-4-16のサブルーチンLDLTの対応する□を参照して下さい。

【後退代入】

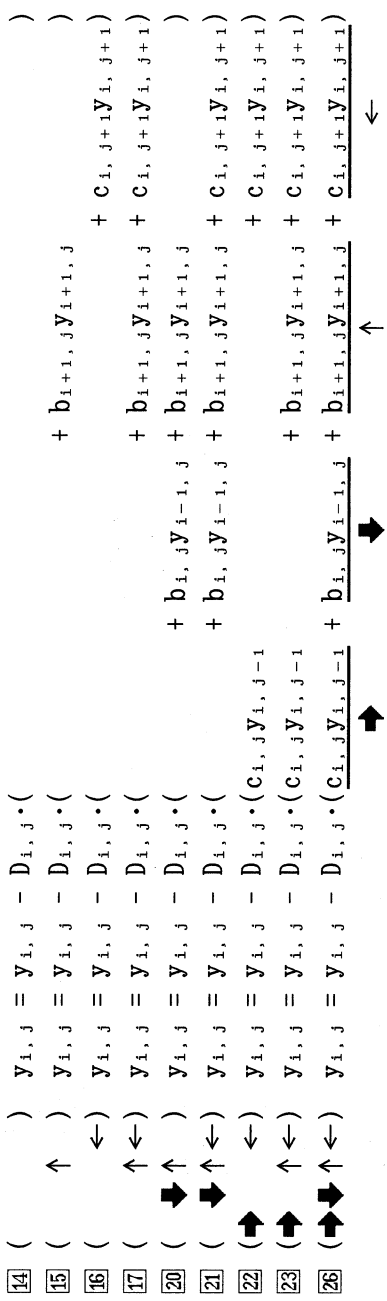


図5-4-15(2)

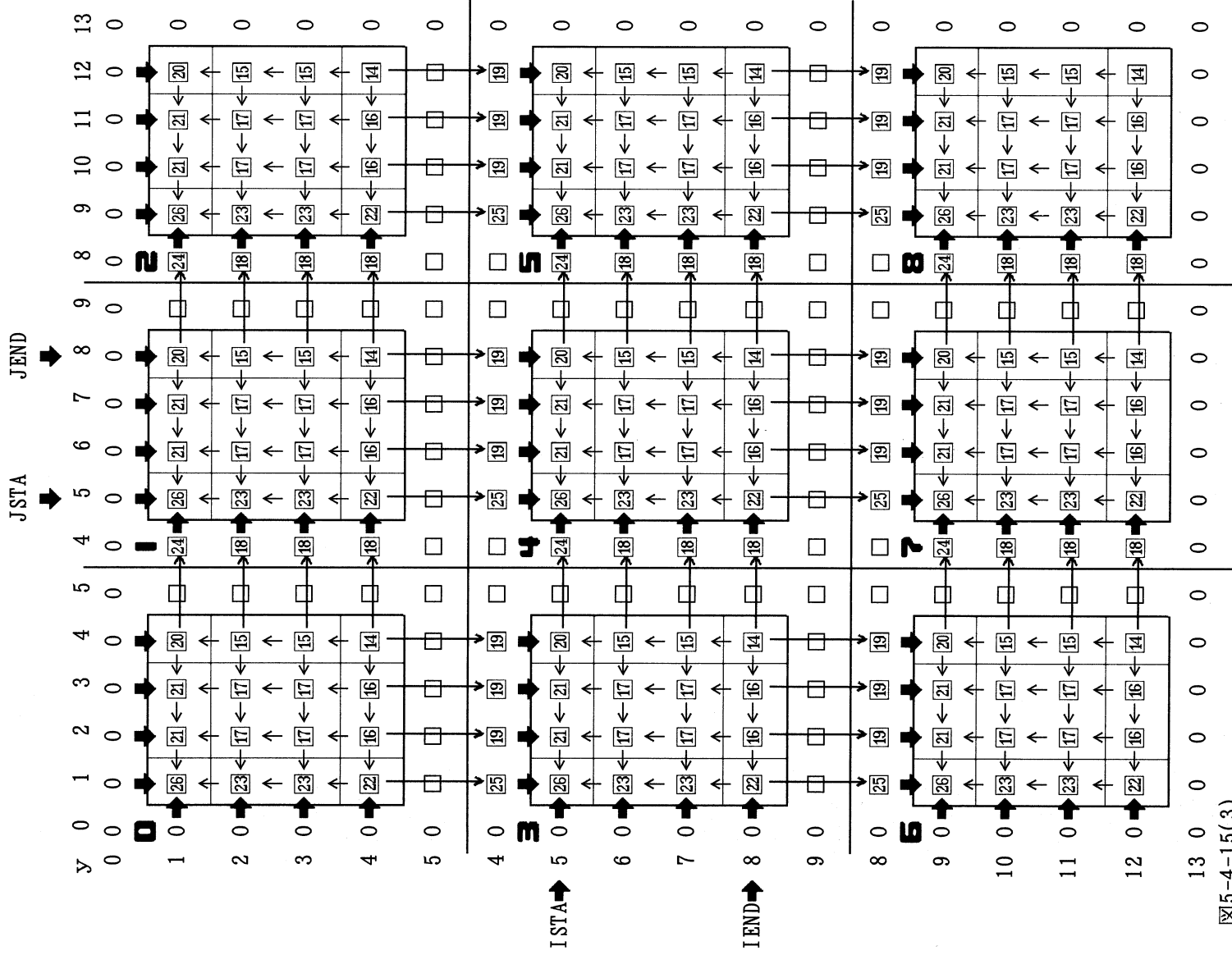


図5-4-15(3)

■ パイプライン法とPBO法の比較

まず収束性について検討します。ICCG法のような反復解法では、収束までの反復回数が少ない方が収束性が高くパフォーマンスが向上します。パイプライン法は単体版のICCG法の前処理をそのまま使用して並列化しているので単体版のICCG法と収束性は変わりませんが、PBO法ではやや不自然な前処理をしているので収束性が若干低下します。またPBO法では分割方法(1次元目と2次元目のプロセス数)によっても収束性が変化します。

次に速度向上率について検討します。パイプライン法は完全な並列化ではないので速度向上率はそれほど高くありませんが、PBO法は完全な並列処理を行っているので連立一次方程式の規模が大きくなればかなり高い速度向上率を得ることができます。



```

MODULE PARA
  INCLUDE 'mpif.h'
  PARAMETER (IPROCS=3, JPROCS=3) [1]
  INTEGER ITABLE(-1:IPROCS, -1:JPROCS)
  INTEGER ITYPE(0:IPROCS*JPROCS-1)
  INTEGER NPROCS, MYRANK, ISTA, IEND, JSTA, JEND
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  INTEGER IUP, IDOWN, JUP, JDOWN
  END

PROGRAM MAIN
  USE PARA
  IMPLICIT REAL*8(A-H, 0-Z)
  PARAMETER (MX=12, MY=12)
  DIMENSION A(0:MX+1, 0:MY+1, 3), B(MX, MY),
  DIMENSION X(0:MX+1, 0:MY+1)
  DIMENSION D(0:MX+1, 0:MY+1), P(0:MX+1, 0:MY+1)
  DIMENSION Q(0:MX+1, 0:MY+1), R(0:MX+1, 0:MY+1)
  DIMENSION BUFS1(MY), BUFR1(MY)
  DIMENSION BUFS2(MY), BUFR2(MY)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(~, NPROCS, IERR)
  CALL MPI_COMM_RANK(~, MYRANK, IERR)
  IF (NPROCS/=IPROCS*JPROCS) THEN
    IF (MYRANK=0) PRINT *, '====ERROR===='
    CALL MPI_FINALIZE(IERR)
  STOP
  ENDF
  DO I=-1, IPROCS
    DO J=-1, JPROCS
      ITABLE(I, J) = MPI_PROC_NULL
    ENDDO
  ENDDO
  IRANK = 0
  DO I=0, IPROCS-1
    DO J=0, JPROCS-1
      ITABLE(I, J) = IRANK
      IF (MYRANK == IRANK) THEN
        MYRANKI = I
        MYRANKJ = J
      ENDIF
    CALL PARA_RANGE
      (1, MX, IPROCS, I, ISTA, IEND)
    CALL PARA_RANGE
      (1, MY, JPROCS, J, JSTA, JEND)
    CALL PARA_TYPE_BLOCK2(0, MX+1, 0, ISTA,
      IEND, JSTA, JEND, MPI_REAL8, ITYPE(IRANK))
    IRANK = IRANK + 1
  ENDDO
  ENDDO

```



```

CALL PARA_RANGE(1, MX, IPROCS, MYRANKI,
  & ISTA, IEND)
CALL PARA_RANGE(1, MY, JPROCS, MYRANKJ,
  & JSTA, JEND) [5]
JUP = ITABLE(MYRANKI, MYRANKJ+1)
JDOWN = ITABLE(MYRANKI, MYRANKJ-1)
IUP = ITABLE(MYRANKI+1, MYRANKJ)
IDOWN = ITABLE(MYRANKI-1, MYRANKJ)
:
  配列AとBに値を設定します。
  EPS=1.0D-7
  ITR=1000
  CALL ICCG(A, MX, MY, B, EPS, ITR, X, D, P, Q, R, IER,
  & BUFS1, BUFR1, BUFS2, BUFR2)
:
  CALL MPI_FINALIZE(IERR)
  END

SUBROUTINE ICCG(A, MX, MY, B, EPS, ITR, X, D, P, Q, R,
  & IER, BUFS1, BUFR1, BUFS2, BUFR2)
  USE PARA
  IMPLICIT REAL*8(A-H, 0-Z)
  DIMENSION A(0:MX+1, 0:MY+1, 3), B(MX, MY)
  DIMENSION X(0:MX+1, 0:MY+1)
  DIMENSION D(0:MX+1, 0:MY+1), P(0:MX+1, 0:MY+1)
  DIMENSION Q(0:MX+1, 0:MY+1), R(0:MX+1, 0:MY+1)
  DIMENSION BUFS1(MY), BUFR1(MY)
  DIMENSION BUFS2(MY), BUFR2(MY)
  DIMENSION WORKS(2), WORKR(2)
  DO J = 0, MY+1
    DO I = 0, MX+1
      D(I, J) = 0.0D0
      X(I, J) = 0.0D0
      P(I, J) = 0.0D0
      Q(I, J) = 0.0D0
      R(I, J) = 0.0D0
    ENDDO
  ENDDO
  CALL DECOMP(A, D, MX, MY, BUFS1, BUFR1)
  CALL AXSUB(A, X, Q, MX, MY,
  & BUFS1, BUFR1, BUFS2, BUFR2)
  DO J = JSTA, JEND
    DO I = ISTA, IEND
      R(I, J) = B(I, J) - Q(I, J)
    ENDDO
  ENDDO
  CALL LDLT(A, D, R, P, MX, MY, BUFS1, BUFR1)

```



⑧ ←

```

CC1 = 0.0D0
DO J = JSTA, JEND
  DO I = ISTA, IEND
    CC1 = CC1 + R(I, J) * P(I, J)
  ENDDO
ENDDO
CALL MPI_ALLREDUCE(CC1, C1, 1, MPI_REAL8,
  & MPI_SUM, MPI_COMM_WORLD, IERR) [6]
&
DO K = 1, ITR  収束反復ループ
  CALL AXSUB(A, P, Q, MX, MY,
    &   BUFS1, BUFR1, BUFS2, BUFR2) [7]
  CC2 = 0.0D0
  DO J = JSTA, JEND
    DO I = ISTA, IEND
      CC2 = CC2 + P(I, J) * Q(I, J)
    ENDDO
  ENDDO
  CALL MPI_ALLREDUCE(CC2, C2, 1, MPI_REAL8, [7]
    & MPI_SUM, MPI_COMM_WORLD, IERR) [7]
  ALPHA = C1 / C2
  X1 = 0.0D0
  X2 = 0.0D0
  DO J = JSTA, JEND
    DO I = ISTA, IEND
      Y = X(I, J)
      X(I, J) = X(I, J) + ALPHA * P(I, J) [10]
      R(I, J) = R(I, J) - ALPHA * Q(I, J) [11]
      X1 = X1 + Y * Y
      X2 = X2 + (X(I, J) - Y) ** 2
    ENDDO
  ENDDO
  WORKS(1) = X1
  WORKS(2) = X2
  CALL MPI_ALLREDUCE(WORKS, WORKR, 2, [8]
    & MPI_REAL8, MPI_SUM, MPI_COMM_WORLD, IERR)
  X1 = WORKR(1)
  X2 = WORKR(2)
  IF (X1 /= 0.0) THEN
    RES = DSQRT(X2/X1)
    IF (RES <= EPS) THEN
      ITR = K
      IER = 0
      GOTO 999
    ENDIF
  ENDIF
  CALL LDLT(A, D, R, Q, MX, MY, BUFS1, BUFR1) [13]

```

→ ⑨

⑩ →

```

CC3 = 0.0D0
DO J = JSTA, JEND
  DO I = ISTA, IEND
    CC3 = CC3 + R(I, J) * Q(I, J)
  ENDDO
ENDDO
CALL MPI_ALLREDUCE(CC3, C3, 1, MPI_REAL8, [9]
  & MPI_SUM, MPI_COMM_WORLD, IERR) [9]
&
BETA = C3/C1
C1 = C3
DO J = JSTA, JEND
  DO I = ISTA, IEND
    P(I, J) = Q(I, J) + BETA * P(I, J)
  ENDDO
ENDDO
IER = 1
999 CONTINUE
DO IRANK=0, NPROCS-1
  CALL MPI_BCAST(X, 1, ITYPE(IRANK), IRANK,
  & MPI_COMM_WORLD, IERR) [10]
  ENDDO
  EPS = RES
  END

```

SUBROUTINE DECOMP(A, D, MX, MY, BUFS1, BUFR1)

```

USE PARA
IMPLICIT REAL*8(A-H, O-Z)
DIMENSION A(0:MX+1, 0:MY+1, 3)
DIMENSION D(0:MX+1, 0:MY+1)
DIMENSION BUFS1(MY), BUFR1(MY)
I = ISTA
J = JSTA
D(I, J) = 1.0D0 / A(I, J, 3) [1]
BUFS1(J) = D(I, J)
DO I = ISTA+1, IEND-1
  D(I, J) = 1.0D0 / ( A(I, J, 3)
    & -D(I-1, J ) * A(I , J , 2) ** 2 ) [2]
  ENDDO
DO J = JSTA+1, JEND-1
  I = ISTA
  D(I, J) = 1.0D0 / ( A(I, J, 3)
    & -D(I , J-1) * A(I , J , 1) ** 2 ) [3]
  BUFS1(J) = D(I, J)
  DO I = ISTA+1, IEND-1
    D(I, J) = 1.0D0 / ( A(I, J, 3)
      & -D(I , J-1) * A(I , J , 1) ** 2
      & -D(I-1, J ) * A(I , J , 2) ** 2 ) [4]
    ENDDO
  ENDDO

```

⑪ ↓

①

```

CALL MPI_ISEND(D(ISTA, JSTA), IEND-ISTA, [5]
& MPI_REAL8, JDOWN, 1,
& MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(BUFS1(JSTA), JEND-JSTA, [6]
& MPI_REAL8, IDOWN, 1,
& MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_IRECV(D(ISTA, JEND+1), IEND-ISTA, [5]
& MPI_REAL8, JUP, 1,
& MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_IRECV(BUFR1(JSTA), JEND-JSTA, [6]
& MPI_REAL8, IUP, 1,
& MPI_COMM_WORLD, IRECV2, IERR)
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
DO J=JSTA, JEND-1
D(IEND+1, J) = BUFR1(J)
ENDDO
I = IEND
J = JSTA
D(I, J) = 1.0D0/( A(I, J, 3)
& -D(I-1, J ) * A(I , J , 2) ** 2
& -D(I+1, J ) * A(I+1, J , 2) ** 2 )
DO J = JSTA+1, JEND-1
D(I, J) = 1.0D0/( A(I, J, 3)
& -D(I, J-1) * A(I, J, 1) ** 2
& -D(I-1, J ) * A(I , J , 2) ** 2
& -D(I+1, J ) * A(I+1, J , 2) ** 2 )
ENDDO
I = ISTA
J = JEND
D(I, J) = 1.0D0/( A(I, J, 3)
& -D(I, J-1) * A(I, J, 1) ** 2
& -D(I, J+1) * A(I, J+1, 1) ** 2 )
DO I = ISTA+1, IEND-1
D(I, J) = 1.0D0/( A(I, J, 3)
& -D(I, J-1) * A(I, J, 1) ** 2
& -D(I-1, J ) * A(I, J, 2) ** 2
& -D(I, J+1) * A(I, J+1, 1) ** 2 )
ENDDO
CALL MPI_ISEND(D(IEND, JSTA), 1, MPI_REAL8, [1]
& JDOWN, 1, MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(D(ISTA, JEND), 1, MPI_REAL8, [2]
& IDOWN, 1, MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_IRECV(D(IEND, JEND+1), 1, MPI_REAL8, [1]
& JUP, 1, MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_IRECV(D(IEND+1, JEND), 1, MPI_REAL8, [2]
& IUP, 1, MPI_COMM_WORLD, IRECV2, IERR)

```

②

③

```

CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
I = IEND
J = JEND
D(I, J) = 1.0D0/( A(I, J, 3)
& -D(I, J-1) * A(I, J, 1) ** 2
& -D(I-1, J ) * A(I, J, 2) ** 2
& -D(I+1, J ) * A(I+1, J, 2) ** 2
& -D(I, J+1) * A(I, J+1, 1) ** 2 )
END
SUBROUTINE AXSUB(A, X, Y, MX, MY,
& BUFS1, BUFR1, BUFS2, BUFR2)
USE PARA
IMPLICIT REAL*8(A-H, O-Z)
DIMENSION A(0:MX+1, 0:MY+1, 3)
DIMENSION X(0:MX+1, 0:MY+1), Y(0:MX+1, 0:MY+1)
DIMENSION BUFS1(MY), BUFR1(MY)
DIMENSION BUFS2(MY), BUFR2(MY)
DO J=JSTA, JEND
BUFS1(J) = X(ISTA, J)
BUFS2(J) = X(IEND, J)
ENDDO
CALL MPI_ISEND(X(ISTA, JSTA), IEND-ISTA+1,
& MPI_REAL8, JDOWN, 1,
& MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(X(ISTA, JEND), IEND-ISTA+1,
& MPI_REAL8, JUP, 1,
& MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_ISEND(BUFS1(JSTA), JEND-JSTA+1,
& MPI_REAL8, IDOWN, 1,
& MPI_COMM_WORLD, ISEND3, IERR)
CALL MPI_ISEND(BUFS2(JSTA), JEND-JSTA+1,
& MPI_REAL8, IUP, 1,
& MPI_COMM_WORLD, ISEND4, IERR)
CALL MPI_IRECV(X(ISTA, JEND+1), IEND-ISTA+1,
& MPI_REAL8, JUP, 1,
& MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_IRECV(X(ISTA, JSTA-1), IEND-ISTA+1,
& MPI_REAL8, JDOWN, 1,
& MPI_COMM_WORLD, IRECV2, IERR)
CALL MPI_IRECV(BUFR1(JSTA), JEND-JSTA+1,
& MPI_REAL8, IUP, 1,
& MPI_COMM_WORLD, IRECV3, IERR)
CALL MPI_IRECV(BUFR2(JSTA), JEND-JSTA+1,
& MPI_REAL8, IDOWN, 1,
& MPI_COMM_WORLD, IRECV4, IERR)

```

④

⤴

```

CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(ISEND3, ISTATUS, IERR)
CALL MPI_WAIT(ISEND4, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV3, ISTATUS, IERR)
CALL MPI_WAIT(IRECV4, ISTATUS, IERR)
DO J=JSTA, JEND
  X(IEND+1, J) = BUFR1(J)
  X(ISTA-1, J) = BUFR2(J)
ENDDO
DO J = JSTA, JEND
  Y(I, J) = A(I, J, 1)*X(I, J-1)
  + A(I, J, 2)*X(I-1, J)
  + A(I, J, 3)*X(I, J)
  + A(I+1, J, 2)*X(I+1, J)
  + A(I, J+1, 1)*X(I, J+1)
ENDDO
ENDDO
END

```

```

SUBROUTINE LDLT(A,D,X,Y,MX,MY,BUFS1,BUFR1)
USE PARA
IMPLICIT REAL*8(A-H,0-Z)
DIMENSION A(0:MX+1,0:MY+1,3)
DIMENSION D(0:MX+1,0:MY+1)
DIMENSION X(0:MX+1,0:MY+1),Y(0:MX+1,0:MY+1)
DIMENSION BUFS1(MY),BUFR1(MY)

```

```

I = ISTA
J = JSTA
Y(I, J) = D(I, J)*X(I, J)
BUFS1(J) = Y(I, J)
DO I = ISTA+1, IEND-1
  Y(I, J) = D(I, J)*( X(I, J)
  -A(I, J, 2)*Y(I-1, J ) )
ENDDO
DO J = JSTA+1, JEND-1
  I = ISTA
  Y(I, J) = D(I, J)*( X(I, J)
  -A(I, J, 1)*Y(I, J-1) )
BUFS1(J) = Y(I, J)
DO I = ISTA+1, IEND-1
  Y(I, J) = D(I, J)*( X(I, J)
  -A(I, J, 1)*Y(I, J-1)
  -A(I, J, 2)*Y(I-1, J ) )
ENDDO
ENDDO

```

⤵

⤵

```

CALL MPI_ISEND(Y(ISTA, JSTA), IEND-ISTA,
& MPI_REAL8, JDOWN, 1,
& MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(BUFS1(JSTA), JEND-JSTA,
& MPI_REAL8, IDOWN, 1,
& MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_IRECV(Y(ISTA, JEND+1), IEND-ISTA,
& MPI_REAL8, JUP, 1,
& MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_IRECV(BUFR1(JSTA), JEND-JSTA,
& MPI_REAL8, IUP, 1,
& MPI_COMM_WORLD, IRECV2, IERR)
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
DO J=JSTA, JEND-1
  Y(IEND+1, J) = BUFR1(J)
ENDDO
I = IEND
J = JSTA
Y(I, J) = D(I, J)*( X(I, J)
& -A(I, J, 2)*Y(I-1, J )
& -A(I+1, J, 2)*Y(I+1, J ) )
DO J = JSTA+1, JEND-1
  Y(I, J) = D(I, J)*( X(I, J)
& -A(I, J, 1)*Y(I, J-1)
& -A(I, J, 2)*Y(I-1, J )
& -A(I+1, J, 2)*Y(I+1, J ) )
ENDDO
I = ISTA
J = JEND
Y(I, J) = D(I, J)*( X(I, J)
& -A(I, J, 1)*Y(I, J-1)
& -A(I, J+1, 1)*Y(I, J+1) )
DO I = ISTA+1, IEND-1
  Y(I, J) = D(I, J)*( X(I, J)
& -A(I, J, 1)*Y(I, J-1)
& -A(I, J, 2)*Y(I-1, J )
& -A(I, J+1, 1)*Y(I, J+1) )
ENDDO
CALL MPI_ISEND(Y(IEND, JSTA), 1, MPI_REAL8,
& JDOWN, 1, MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(Y(ISTA, JEND), 1, MPI_REAL8,
& IDOWN, 1, MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_IRECV(Y(IEND, JEND+1), 1, MPI_REAL8,
& JUP, 1, MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_IRECV(Y(IEND+1, JEND), 1, MPI_REAL8,
& IUP, 1, MPI_COMM_WORLD, IRECV2, IERR)

```

⤴



④

```

CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
I = IEND
J = JEND
Y(I,J) = D(I,J)*( X(I,J)
& -A(I , J , 1)*Y(I , J-1)
& -A(I , J , 2)*Y(I-1,J )
& -A(I+1,J , 2)*Y(I+1,J )
& -A(I , J+1,1)*Y(I , J+1) )
I = IEND
J = JEND
BUFS1(J) = Y(I,J)
D0 I = IEND-1, ISTA+1, -1
Y(I,J) = Y(I,J)-D(I,J)*(
& +A(I+1,J , 2)*Y(I+1,J ) )
ENDDO
D0 J = JEND-1, JSTA+1, -1
I = IEND
Y(I,J) = Y(I,J)-D(I,J)*(
& +A(I , J+1,1)*Y(I , J+1) )
BUFS1(J) = Y(I,J)
D0 I = IEND-1, ISTA+1, -1
Y(I,J) = Y(I,J)-D(I,J)*(
& +A(I+1,J , 2)*Y(I+1,J )
& +A(I , J+1,1)*Y(I , J+1) )
ENDDO
CALL MPI_ISEND(Y(ISTA+1,JEND), IEND-ISTA,
& MPI_REAL8, JUP, 1,
& MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(BUFS1(JSTA+1) , JEND-JSTA,
& MPI_REAL8, IUP, 1,
& MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_Irecv(Y(ISTA+1, JSTA-1), IEND-ISTA,
& MPI_REAL8, JDOWN, 1,
& MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_Irecv(BUFR1(JSTA+1) , JEND-JSTA,
& MPI_REAL8, IDOWN, 1,
& MPI_COMM_WORLD, IRECV2, IERR)
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
D0 J=JSTA+1, JEND
Y(ISTA-1, J) = BUFR1(J)
ENDDO

```

①

①

```

I = ISTA
J = JEND
Y(I,J) = Y(I,J)-D(I,J)*(
& +A(I , J , 2)*Y(I-1,J )
& +A(I+1,J , 2)*Y(I+1,J ) )
D0 J = JEND-1, JSTA+1, -1
Y(I,J) = Y(I,J)-D(I,J)*(
& +A(I , J , 2)*Y(I-1,J )
& +A(I+1,J , 2)*Y(I+1,J )
& +A(I , J+1,1)*Y(I , J+1) )
ENDDO
I = IEND
J = JSTA
Y(I,J) = Y(I,J)-D(I,J)*(
& A(I , J , 1)*Y(I , J-1)
& +A(I , J+1,1)*Y(I , J+1) )
D0 I = IEND-1, ISTA+1, -1
Y(I,J) = Y(I,J)-D(I,J)*(
& A(I , J , 1)*Y(I , J-1)
& +A(I+1,J , 2)*Y(I+1,J )
& +A(I , J+1,1)*Y(I , J+1) )
ENDDO
CALL MPI_ISEND(Y(ISTA,JEND), 1, MPI_REAL8,
& JUP, 1, MPI_COMM_WORLD, ISEND1, IERR)
CALL MPI_ISEND(Y(IEND, JSTA), 1, MPI_REAL8,
& IUP, 1, MPI_COMM_WORLD, ISEND2, IERR)
CALL MPI_Irecv(Y(ISTA, JSTA-1), 1, MPI_REAL8,
& JDOWN, 1, MPI_COMM_WORLD, IRECV1, IERR)
CALL MPI_Irecv(Y(ISTA-1, JSTA), 1, MPI_REAL8,
& IDOWN, 1, MPI_COMM_WORLD, IRECV2, IERR)
CALL MPI_WAIT(ISEND1, ISTATUS, IERR)
CALL MPI_WAIT(ISEND2, ISTATUS, IERR)
CALL MPI_WAIT(IRECV1, ISTATUS, IERR)
CALL MPI_WAIT(IRECV2, ISTATUS, IERR)
I = ISTA
J = JSTA
Y(I,J) = Y(I,J)-D(I,J)*(
& A(I , J , 1)*Y(I , J-1)
& +A(I , J , 2)*Y(I-1,J )
& +A(I+1,J , 2)*Y(I+1,J )
& +A(I , J+1,1)*Y(I , J+1) )
END

```

図5-4-16

## 5-5 マルチフロントナル法

不規則疎行列の連立一次方程式の直接法として、一般にスカイライン法が使用されます。図5-5-1(1)のよ  
うな不規則疎行列(本例は対称行列を想定)の場合、スカイライン法では着色した部分のみをメモリー上に保  
持するので、行列全体を密行列として取り扱うよりも計算量が低減され、メモリーも節約できます。

ところが、図5-5-1(2)のようにバンド幅の広い不規則疎行列の場合、対角要素から遠く離れた要素が一つ  
でもあると、その要素から対角要素までの全ての要素を、途中が0でも全て保持する必要があるため、計算  
量はあまり少なくならず、メモリーも節約できません。

近年マルチフロントナル法という解法が脚光を浴びており、市販の構造解析などのパッケージにも採用され  
ています。マルチフロントナル法は、本例のようにバンド幅の広い不規則疎行列の場合、従来のスカイライン法と比  
(これ以外に作業域が必要)、本例のようにバンド幅の広い不規則疎行列の場合、従来のスカイライン法と比  
べて計算時間が短縮でき、メモリーも少なくとも紹介します。

マルチフロントナル法の並列版は、第6章で紹介する数値計算ライブラリーの中に含まれている可能性があ  
りますが、未調査です。

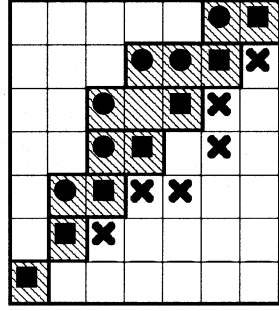


図5-5-1(1) バンド幅が狭い  
スカイライン法

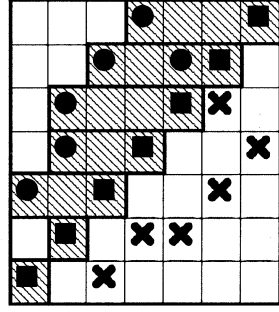


図5-5-1(2) バンド幅が広い  
スカイライン法

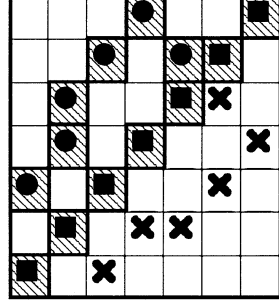


図5-5-1(3) バンド幅が広い  
マルチフロントナル法

## 5-6 SOR法

### 5-6-1 SOR法と並列性

差分法で使用されるガウスザイデル法やその改良版であるSOR法は、そのままでは並列性がありませんが、計算順序を変更することによって並列化することができます。まず、ガウスザイデル法はなぜ並列性がないのか、そして並列性を持たせるにはどのような方法によればよいかについて図5-6-1で説明します。

● 図5-6-1(次ページ)の(1)に示す、未知数が $x_1 \sim x_4$ の1次元の計算領域に対し、(2)に示す1次元ラプラス方程式を適用し、(3)のように差分化します。(3)を各要素に対して書き下すと(4)になり、これを上から順に計算したとすると、例えば1番目の式で求められた $x_1$ が2番目の式の右辺で参照されています。このような依存関係を $\searrow$ で表わすことにします。すると(4)の各式には、 $i$ 番目の式の計算が終了しないと $i+1$ 番目の式の計算を始めることができないという依存関係があり、並列性のないことが分かります。なお、(4)の順序で計算を行うのが、(5)に示す通常のガウスザイデル法です。

● (4)を(6)のように書き表してみます。ここで例えば $x_1$ の計算前の値を①、計算後の値を②と表し、境界の定数を「境」と表します。例えば②の計算に①と③が使用されるので、これを $\checkmark$ と $\rightarrow$ で表します。例えば②のようにランク0と1の2つのプロセスで並列する場合、②から③へプロセス間にもたがることがあり、②の計算が終わらないと③を計算できないため、並列に実行することはできないことが分かります。

● (4)では上から順に計算しましたが、計算順序は必ずしもこの順である必要はありません。そこで、(4)の2式目と3式目を入れ替えた(7)を上から順に計算してみます。すると $\searrow$ の矢印の状態は先ほどの(4)と変わりました。しかしこの図からでは、並列性があるのかどうかはつきりしません。

そこで(7)を(8)で表し、さらに書き方を少し変えた(9)で検討してみます。すると①と③の計算、および②と④の計算を、ランク0と1で並列に行うことができることが分かります。ただし、①と③の計算が終了してから②と④の計算を行う必要があり、またランク0とランク1の間に矢印がまたがっているため、プロセス間で通信が必要となります。通信の部分をより明確にしたのが(10)で、 $\blackrightarrow$ と $\blackleftarrow$ は通信を表します。

なお、(4)と(7)は計算順序が異なるので、計算を1回行っただけでは計算結果は異なります。しかし計算を何度も反復させると次第に解は同じ値に収束していきます。ただし(4)と(7)で収束性は変わります。

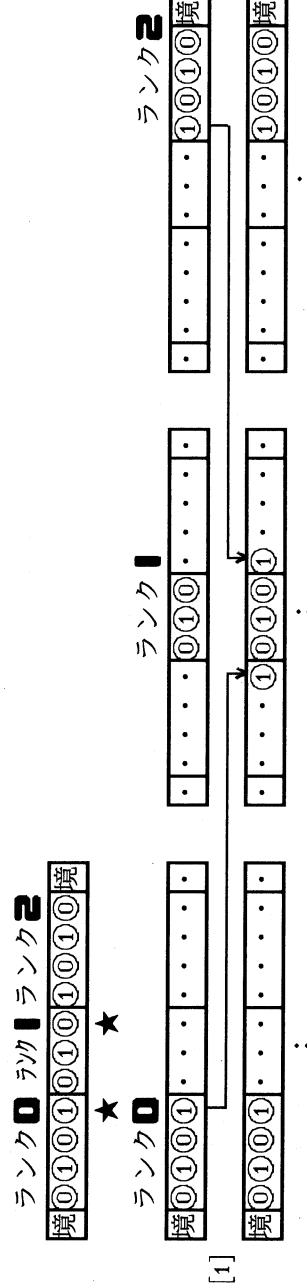
● 以後、(10)を(11)のように表すことにします。ランク0は $x_1$ と $x_2$ を担当し、ランク1は $x_3$ と $x_4$ を担当します。  
 [1] まず、プロセス間の境界のデータ②を $\blackleftarrow$ 方向へ通信します。

[2] 各プロセスは奇数番号の要素①③を計算します。

[3] プロセス間の境界のデータ③を $\blackrightarrow$ 方向へ通信します。

[4] 各プロセスは偶数番号の要素②④を計算します。

● もう少し要素数が多い例を(12)に示します。(12)を上記の方法で3つのプロセスで並列に計算したときの動作を(13)に示します。以後の説明の都合上、要素の番号は①②③④でなく①①①①と表します。①と①は計算順序を意味します。[2]ではまず①の要素を計算し、[4]では①の要素を計算します。この方法を本書ではマルチカラー法と呼びます(計算領域が2つの異なる部分に互い違いに分かれるので)。なお(12)で各プロセス(ただしランク2を除く)の右側の境界の要素(★)は全て①になっています。理由は、右側の境界に①と①が混合していると、下図のようにプロセスによって通信の方向が反対になってしまい、並列化しにくいからです。



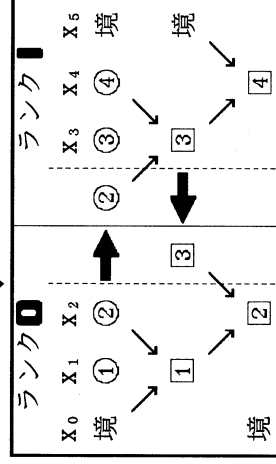
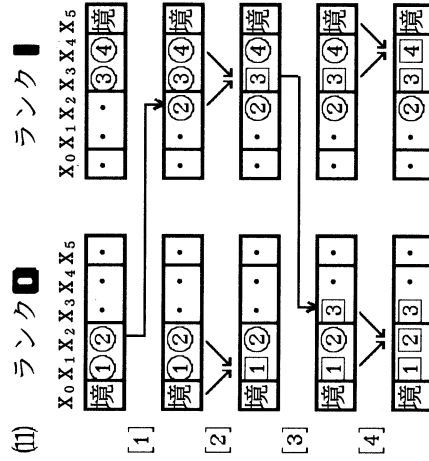
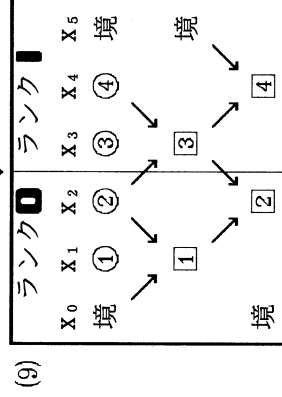
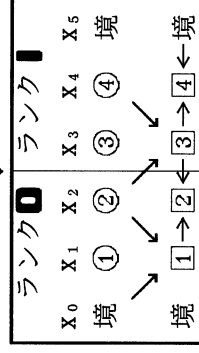
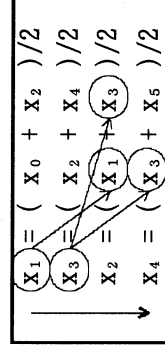
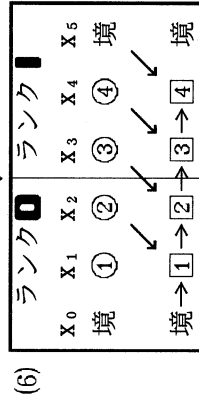
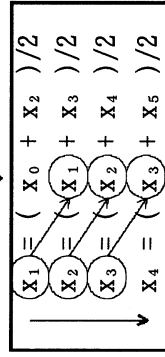
$X_0 X_1 X_2 X_3 X_4 X_5$

(1) (計算領域) 境①②③④境

(2) (微分方程式)  $d^2 x/dy^2 = 0$

(3) (差分式)  $X_i = (X_{i-1} + X_{i+1})/2$

```
do k=1,10000
do i=1,4
  Xi=(Xi-1+Xi+1)/2
enddo
収束したかチェック
enddo
```



(12) ランク0 ランク1 ランク2

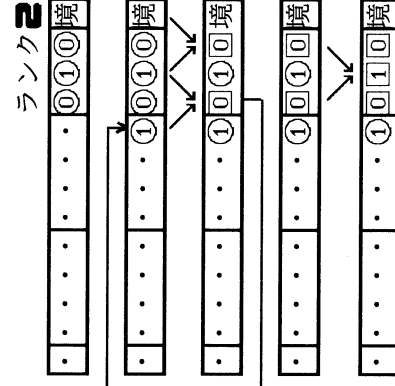
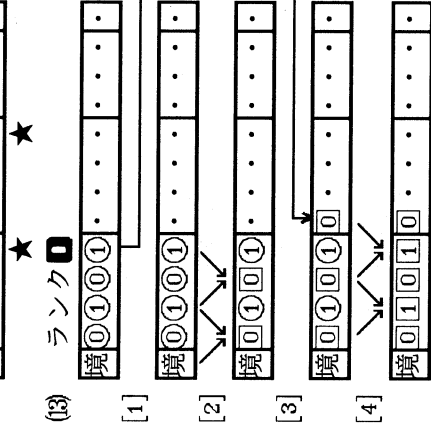


図5-6-1

## 5-6-2 プログラム例 (ゼブラ法)

本節では、2次元のガウスサイズル法やSOR法のプログラムを実際に並列化する方法を説明します。図5-6-2(1)に2次元SOR法のプログラムを示します。図5-6-3(1)が計算領域で、「境」は境界部分の要素を示します。計算は矢印の順に行われ、●は計算が終了している要素、○は現在計算している要素、◎は現在計算している要素、○はこれから計算する要素を示します。また◎の要素の計算には+で囲んだ部分の要素を使用します(5点差分)。

このプログラムを、まず前述のマルチカラー法の単体版に変更します。計算領域が1次元の場合、各要素を図5-6-3(2)の上段に示すように◎と①に分割しましたが、2次元の場合はこれをI方向に引き伸ばしたと考え、図5-6-3(2)の下段のように分割します。◎と①が列単位に互い違いになっているので、本書ではこれをゼブラ法と呼びます。計算順序はまず図5-6-3(2)で◎を計算し、次に図5-6-3(3)で①を計算します。

ゼブラ法(単体版)のプログラムを図5-6-2(2)に示します。オリジナル版との変更点を下線で示します。  
 [0]のD0ループがJJ=0のときに図5-6-3(2)の◎を、JJ=1のとき図5-6-3(3)の①を計算します。

```

:
PARAMETER (IMIN=1, IMAX=4, JMIN=1, JMAX=13)
DIMENSION X(IMIN-1:IMAX+1, JMIN-1:JMAX+1)
:
D0 K=1, 10000 (収束反復)
ERR = 0.0
D0 J=JMIN, JMAX
D0 I=IMIN, IMAX
TEMP = 0.25*
(X(I+1, J) + X(I, J+1)
+ X(I, J-1) + X(I-1, J)) - X(I, J)
X(I, J) = X(I, J) + OMEGA * TEMP
ERR = MAX(ERR, ABS(TEMP))
ENDDO
ENDDO
IF (ERR <= EPS) EXIT (収束判定)
ENDDO
:
    
```

図5-6-2(1) オリジナル

```

:
PARAMETER (IMIN=1, IMAX=4, JMIN=1, JMAX=13)
DIMENSION X(IMIN-1:IMAX+1, JMIN-1:JMAX+1)
:
D0 K=1, 10000 (収束反復)
ERR = 0.0
D0 JJ=0, 1
D0 J=JMIN+JJ, JMAX, 2
D0 I=IMIN, IMAX
TEMP = 0.25*
(X(I+1, J) + X(I, J+1)
+ X(I, J-1) + X(I-1, J)) - X(I, J)
X(I, J) = X(I, J) + OMEGA * TEMP
ERR = MAX(ERR, ABS(TEMP))
ENDDO
ENDDO
IF (ERR <= EPS) EXIT (収束判定)
ENDDO
:
    
```

図5-6-2(2) ゼブラ法(単体版)

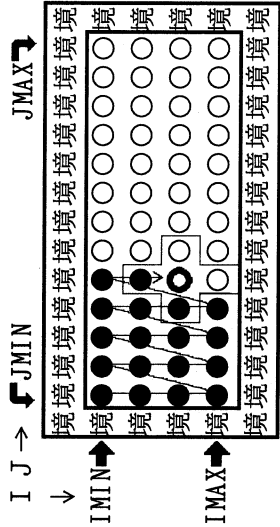


図5-6-3(1) SOR法オリジナル

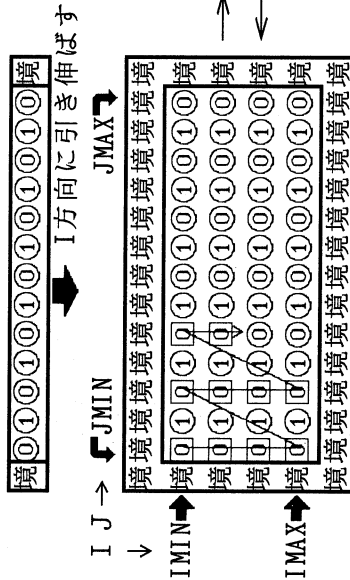


図5-6-3(2) ゼブラ法(単体) ◎の計算

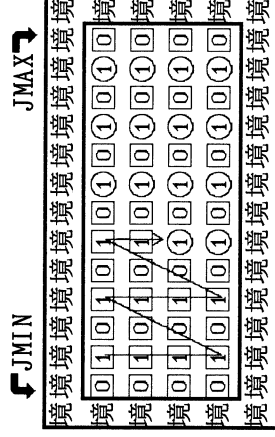


図5-6-3(3) ゼブラ法(単体) ①の計算

次にゼブラ法(単体版)を並列化する方法を説明します。まず図5-6-4に3プロセスで並列化したときの動作を示します。図5-6-1と同様に、分割方法はJ方向でブロック分割します。紙面の関係で「境」の部分は省略し、I方向の要素数は2個とします。ゼブラ法で並列化したプログラムを図5-6-6に示し、説明を以下に示します。この方法は図5-6-10(1)に示すように、2次元の5,9点差分、3次元の7,19点差分の場合に使用できます。

- 図5-6-6の[1]で、プログラムの後半([9]と[10])で行う通信(シフト)の準備をします(4-6-2節参照)。
- [2]～[6]で各プロセスのJ方向の担当範囲を求めます。このとき図5-6-1(2)で説明したように、各プロセスが行う通信方向を一定にするため、各プロセス(ただしランク2を除く)の右側の境界を全て①にする必要があります。この方法を図5-6-5を使用して以下で説明します。
  - ①と②の2列を単位として分割するので、[2]でNBLOCK=2を設定します。
  - [3]で、(2列を1組とした)列数NTOTAL(本例では7)を求めます。
  - [4]で、サブルーチンPARAM\_RANGE(4-5-4節参照)を使用して、1～NTOTALのうちランクがIRANKのプロセスが担当する列番号の下限JSTAXと上限JENDXをブロック分割で求めます。
  - [5]で、ランクがIRANKのプロセスが担当する実際の列番号の下限JJSTA(IRANK)と上限JJEND(IRANK)を求めます。各プロセスは最低2列を担当する必要があるもので、これを[6]でチェックし、2列より少ないプロセスがある場合はエラーで終了します。
  - [7]で改めて、自プロセスが担当するJ方向の下限と上限をJSTA, JENDに設定します。

- [8]のJJが0のとき、[9]で図5-6-4の通信を行い、[11]で各プロセスは②を計算します。次に[8]のJJが1のとき、[10]で通信を行い、[11]で各プロセスは①を計算します。
- [12]で、各プロセスが求めたERRの最大値をMPI\_ALLREDUCEを使用して求めます。

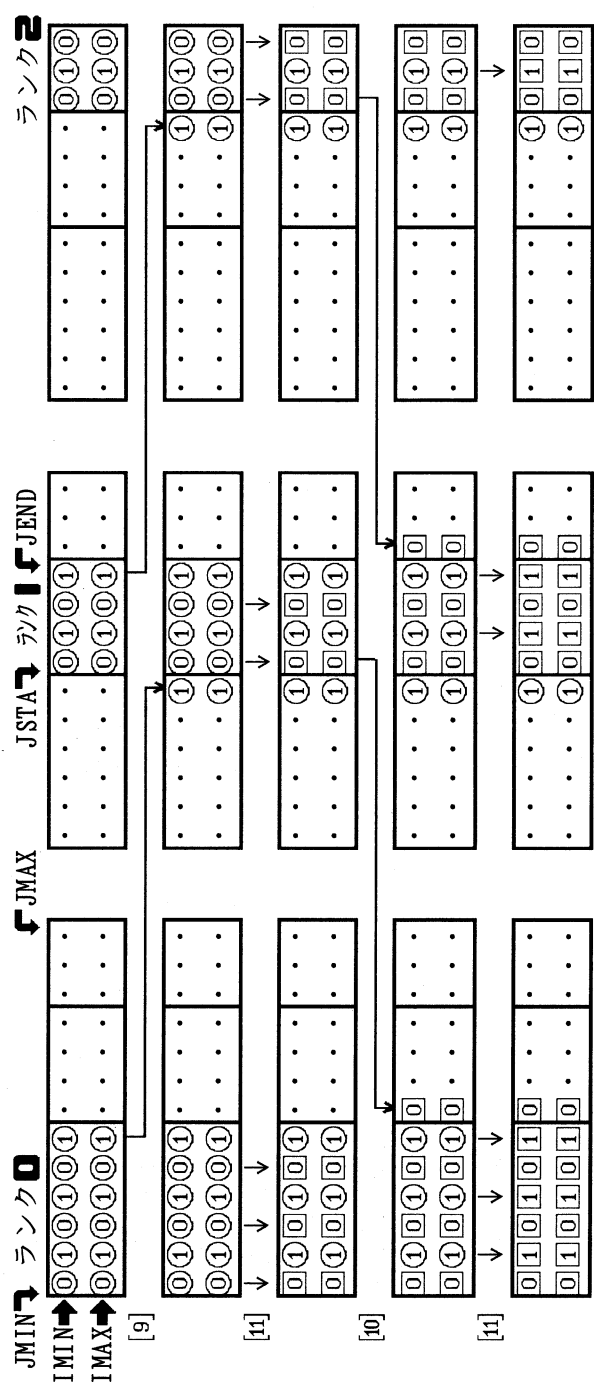
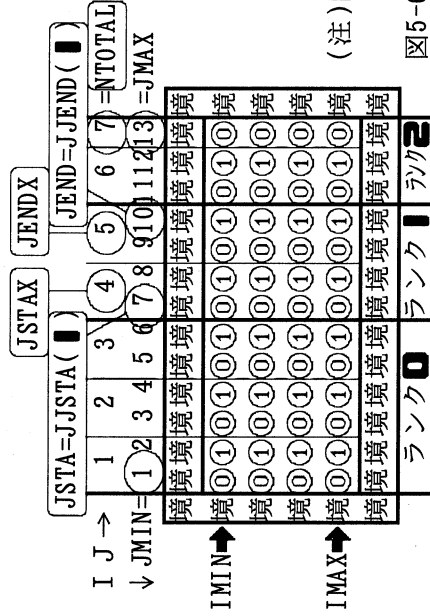


図5-6-4 ゼブラ法(並列版)の動作



(注) 図中のJSTAX, JENDX, JSTA, JENDはランク 1 用です。

図5-6-5 ゼブラ法(並列版)の分割

```

:
INCLUDE 'mpif.h'
PARAMETER (IMIN=1, IMAX=4, JMIN=1, JMAX=13)
DIMENSION X(IMIN-1:IMAX+1, JMIN-1:JMAX+1)
INTEGER ISTATUS(MPI_STATUS_SIZE)
INTEGER, ALLOCATABLE :: JJSTA(:), JJEND(:)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE
& (MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK
& (MPI_COMM_WORLD, MYRANK, IERR)
IUP = MYRANK + 1
IDOWN = MYRANK - 1
IF (MYRANK==NPROCS-1) IUP = MPI_PROC_NULL
IF (MYRANK=0) IDOWN = MPI_PROC_NULL
ALLOCATE(JJSTA(0:NPROCS-1),
& JJEND(0:NPROCS-1))
NBLOCK = 2
NTOTAL = (JMAX-JMIN)/NBLOCK + 1
DO IRANK=0, NPROCS-1
CALL PARA_RANGE(1, NTOTAL, NPROCS,
& IRANK, JSTAX, JENDX)
JJSTA(IRANK) = JMIN+(JSTAX-1)*NBLOCK
JJEND(IRANK)
& = MIN(JMIN+JENDX*NBLOCK-1, JMAX)
IF (JJEND(IRANK)-JJSTA(IRANK)+1
& < NBLOCK) THEN
IF (MYRANK==0) PRINT *, '===ERROR===',
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
ENDDO
JJSTA = JJSTA(MYRANK)
JEND = JJEND(MYRANK)
:

```

図5-6-6 ゼブラ法(並列版)

```

DO K=1, 10000 (収束反復)
ERR = 0.0
DO JJ=0, 1
IF (JJ==0) THEN
CALL MPI_ISEND(X(IMIN-1, JEND),
IMAX-IMIN+3, MPI_REAL, IUP, 1,
MPI_COMM_WORLD, IREQS, IERR)
CALL MPI_IRECV(X(IMIN-1, JSTA-1),
IMAX-IMIN+3, MPI_REAL, IDOWN, 1,
MPI_COMM_WORLD, IREQR, IERR)
ELSE
CALL MPI_ISEND(X(IMIN-1, JSTA),
IMAX-IMIN+3, MPI_REAL, IDOWN, 1,
MPI_COMM_WORLD, IREQS, IERR)
CALL MPI_IRECV(X(IMIN-1, JEND+1),
IMAX-IMIN+3, MPI_REAL, IUP, 1,
MPI_COMM_WORLD, IREQR, IERR)
ENDIF
CALL MPI_WAIT(IREQS, ISTATUS, IERR)
CALL MPI_WAIT(IREQR, ISTATUS, IERR)
DO J=JSTA+J, JEND, 2
DO I=IMIN, IMAX
TEMP = 0.25*
(X(I+1, J) + X(I, J+1)
+ X(I, J-1) + X(I-1, J)) - X(I, J)
X(I, J) = X(I, J) + OMEGA*TEMP
ERR = MAX(ERR, ABS(TEMP))
ENDDO
ENDDO
CALL MPI_ALLREDUCE(ERR, ERR2, 1, MPI_REAL,
MPI_MAX, MPI_COMM_WORLD, IERR)
ERR = ERR2
IF (ERR <= EPS) EXIT (収束判定)
ENDDO
CALL MPI_FINALIZE(IERR)
:

```

## 5-6-3 その他の解法

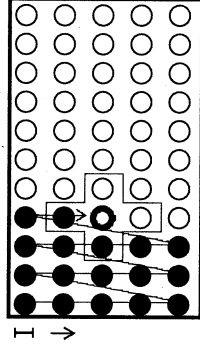
本節ではゼブラ法以外の方法について紹介します。各解法のまとめを図5-6-10に示します。図5-6-10には計算領域が3次元の場合の図を描いてありますが、2次元の場合は一番手前の面ののみを考えて下さい。また、具体的な並列化の方法については5-6-4節で説明します。

### ■ マルチカラー法

● 図5-6-7(1)と(2)に示すオ리지ナルのSOR法は、前述のように並列性がないため、(3)と(4)のように、J方向の計算順序を「①④①④…」と互い違いにして並列性を持たせ、並列化しました(ゼブラ法)。  
 ところでベクトル計算機の場合は、最も内側のループ(本例ではD0 I)がベクトル化の対象になります。(4)の0で囲んだ各要素が独立に計算できればD0 Iがベクトル化できますが、I方向に関しては前述と同じ理由で依存関係があるためこのままではベクトル化できません。そこで(5)のように、I方向の計算順序も「①④①④…」のようにすれば、前述と同じ理由で依存関係がなくなり、ベクトル化できます。(5)のように、IとJの両方向で「①④①④…」の計算順序になっている計算方法を本書では2色のマルチカラー法と呼びます。  
 ところで(5)を並列化した場合、境界の1列のデータ(①または④)を通信しますが、例えば④を通信する場合、④がとびとびに入っていて通信しにくいので、実際には(6)(7)のように4色のマルチカラー法にした方が簡単になります(具体的な通信方法は後述します)。

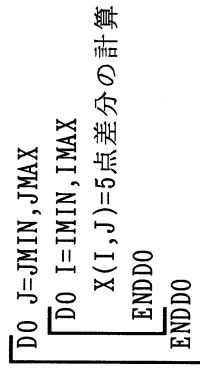
一方スカラ計算機で並列化する場合は(3)(4)のゼブラ法でも良いですが、I方向とJ方向の計算を同じパターンにしたい場合、あるいは(8)のようにI方向とJ方向の両方でブロック分割する場合、マルチカラー法を使用します。ただし、メモリー上でデータの並ぶI方向の要素の計算順序がとびとびになるため、ゼブラ法よりもややキャッシュミスが多く、パフォーマンスが若干低下する可能性があります。

(1) J →

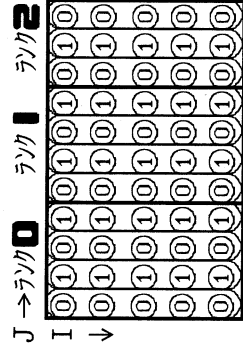


=

(2)

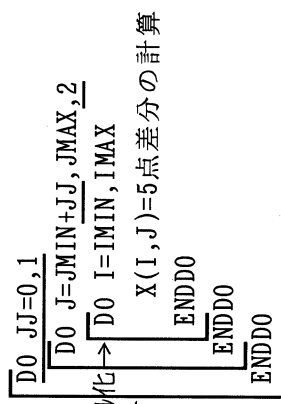


(3) J →

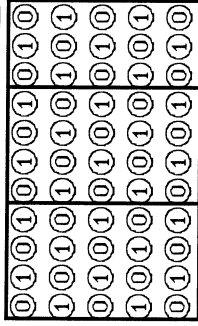


=

(4)

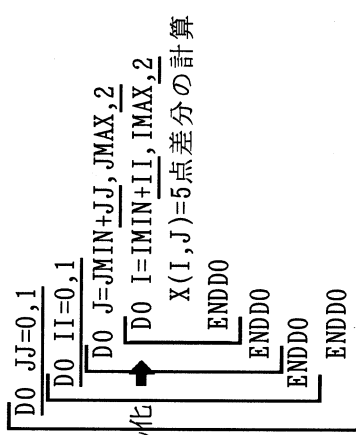


(5) J →



=

(7)



(8) J →

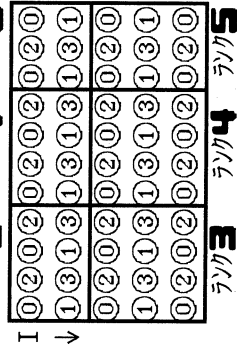


図5-6-7



■ パラレル・ブロック・オーダリング法

1次元のガウスザイデル法で要素数が4個の場合、ゼブラ法やマルチカラー法では図5-6-8(1)の実線のように「①②③④」と分割しました。要素数が多い場合は「①②③④」がそのまま続いたと考え、(2)のように分割しました。

これに対し、要素数が多い場合は「①②③④」のうち③の部分だけを広げ、(5)のように分割する方法を本書ではパラレル・ブロック・オーダリング法と呼びます。(2)と比べ、③の計算がメモリー上で連続するため、キャッシュミスが若干減少すると思われれます。3プロセスで並列化した場合を(3)と(6)に示しますが、パラレル・ブロック・オーダリング法ではプロセス数によって計算パターンが変わるといふ短所があります。

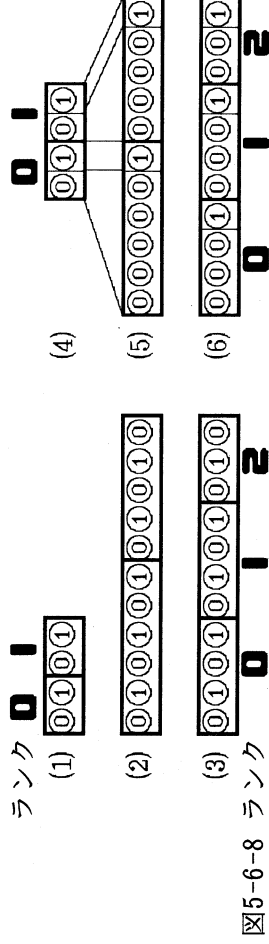


図5-6-8

■ パイプライン法

上記で説明した方法はいずれも、元のガウスザイデル法の計算順序を変更することによって並列性を持たせました。4-6-7節で説明したパイプライン法を利用すると、以下のように、元の計算順序を変えずにとなく並列化できる場合があります(プログラム例は省略します)。

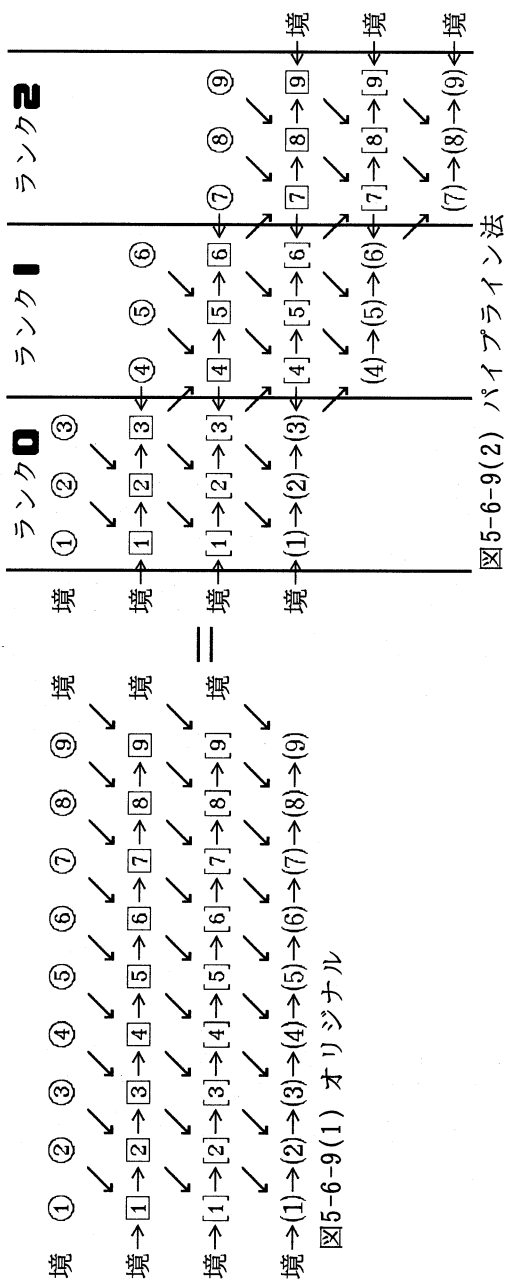


図5-6-9(2) パイプライン法

■ 計算領域(2次元と3次元)の相違

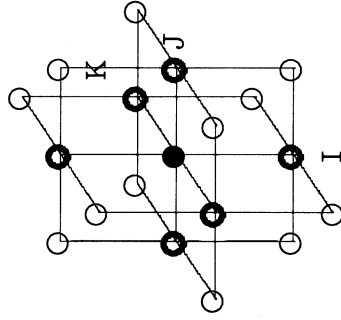
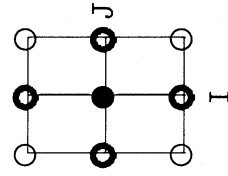
ゼブラ法とパラレル・ブロック・オーダリング法の場合は、図5-6-10の(1)(2)(3)(7)(8)に示すように、2次元の計算領域をK方向に引き伸ばすと3次元になり、計算方法は同じです。マルチカラー法の場合は、図5-6-10の(4)(6)に示すように、2次元と3次元で色の数が変わります。

■ 差分式による相違

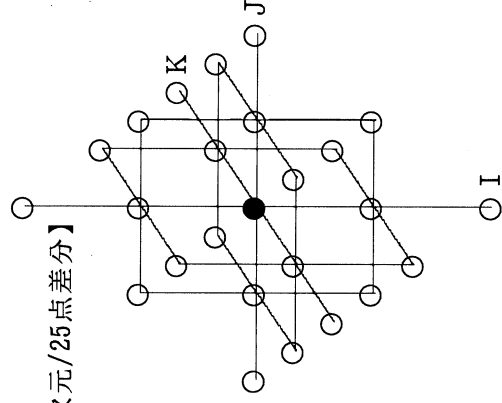
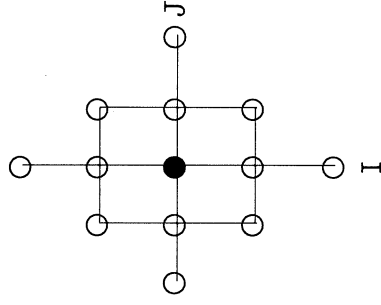
以下の(a)のように、●の計算で1つ離れた値を参照する場合(図5-6-10の左半分)と、(b)のように2つ離れた値も参照する場合(図5-6-10の右半分)で使用する色数が変わり、また並列化したときの通信方法も変わります。具体的な方法については後述します。

$$\begin{array}{l}
 X_{i-2} \quad X_{i-1} \quad X_i \quad X_{i+1} \quad X_{i+2} \\
 (a) \quad \bigcirc - \bullet - \bigcirc - \bigcirc \quad X_i = \quad X_{i-1} + X_{i+1} \\
 (b) \quad \bigcirc - \bigcirc - \bullet - \bigcirc - \bigcirc \quad X_i = \quad X_{i-2} + X_{i-1} + X_{i+1} + X_{i+2}
 \end{array}$$

【2次元/5,9点差分】



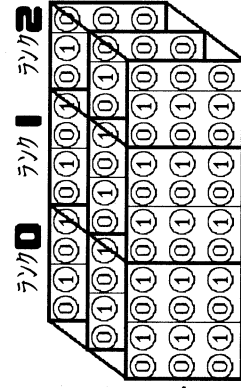
【2次元/13点差分】



【3次元/25点差分】

(1)ゼブラ

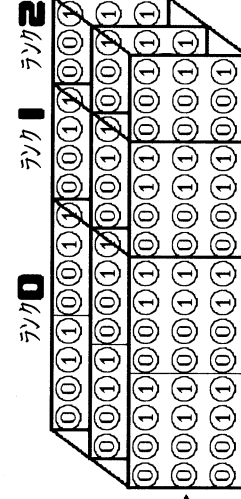
(2色)



↑ KIJ →  
↓

(2)ゼブラ

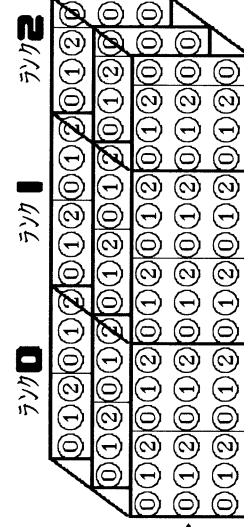
(2色2列)



↑ KIJ →  
↓

(3)ゼブラ

(3色)

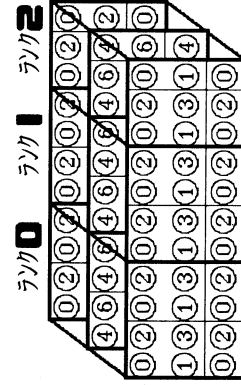


↑ KIJ →  
↓

(4)マルチカラー

(2次元:4色)

(3次元:8色)



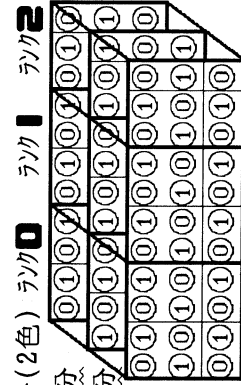
↑ KIJ →  
↓

(5)マルチカラー

(2次元:5点差分)

3次元は7点差分

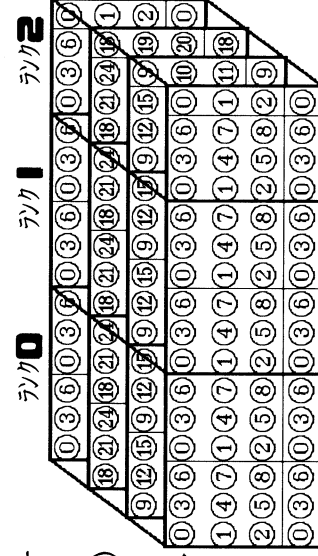
のみ



(6)マルチカラー

(2次元:9色)

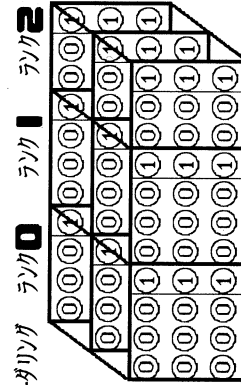
(3次元:27色)



↑ KIJ →  
↓

(7)パレル・ブロック・オーダリング

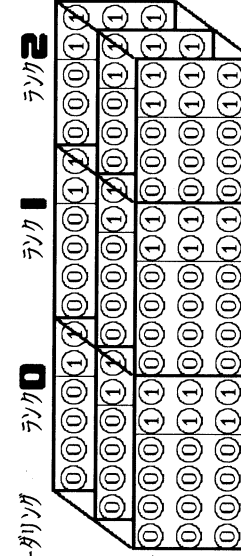
(境界1列)



↑ KIJ →  
↓

(8)パレル・ブロック・オーダリング

(境界2列)



↑ KIJ →  
↓

### 5-6-4 各解法の具体的な並列化方法

本節では、図5-6-10の各方法(5)と図5-6-6で説明した(1)を除く)を並列化した場合のプログラム例とデータの動きについて説明します。

図5-6-6の2次元プログラムのうち、各解法によって異なる部分は[2]~[6]と[8]~[11]なので、プログラム例ではその部分のみを示します。なお[6]のエラーチェック部分と[9][10]の通信ルーチンは省略しました。また3次元プログラムも省略しました。

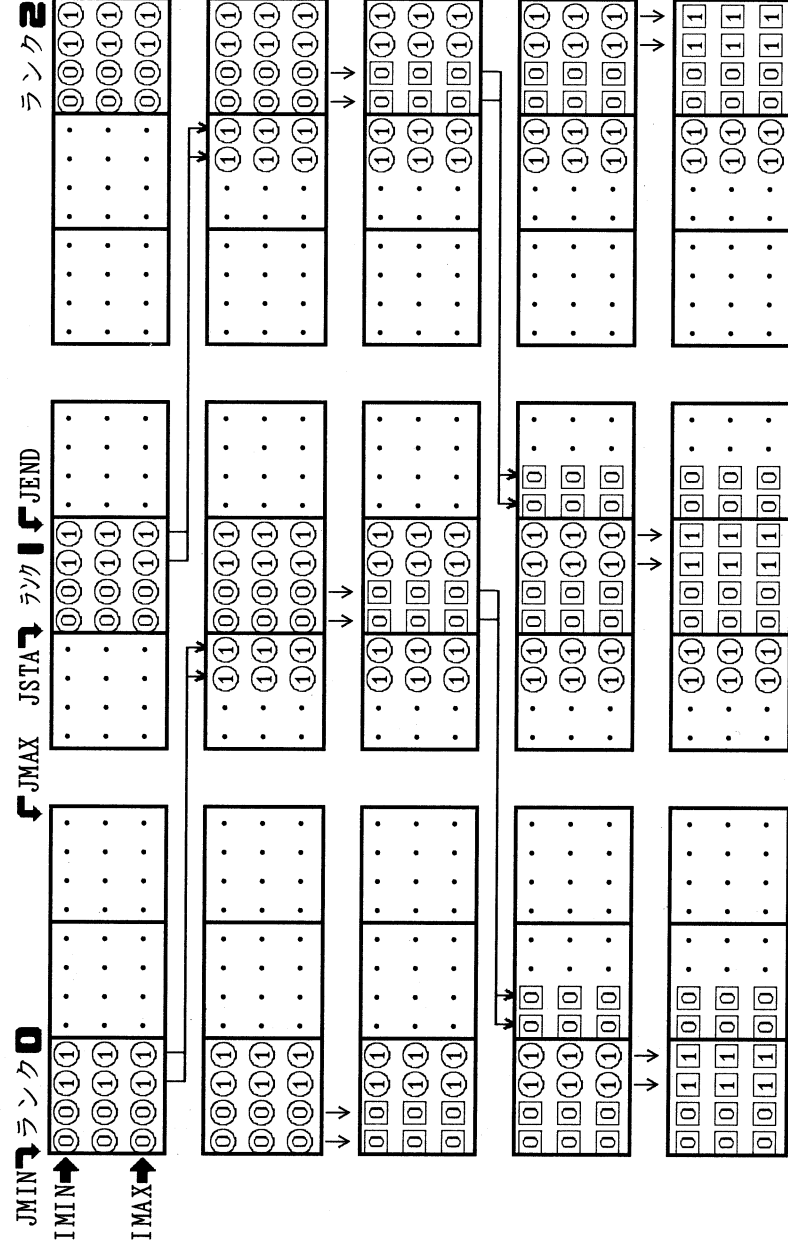
■ ゼブラ (図5-6-10の(2)参照)

- 2次元の13点差分、3次元の25点差分の場合に使用します。
- 各プロセス(ランク2)を除く)が担当する列数は、4列(①②③④)の倍数になっている必要があります。このため[1]で4を指定し、4列を単位としてブロック分割します。
- 各プロセスが担当する列数は最低4列(①②③④)必要なので、これを[2]でチェックします。

```

:
NBLOCK = 4
NTOTAL = (JMAX-JMIN)/NBLOCK + 1
DO IRANK=0,NPROCS-1
CALL PARA_RANGE(1,NTOTAL,NPROCS,IRANK,JSTAX,JENDX)
JJSTA(IRANK) = JMIN+(JSTAX-1)*NBLOCK
JJEND(IRANK) = MIN(JMIN+JENDX*NBLOCK-1,JMAX)
IF (JJEND(IRANK)-JJSTA(IRANK)+1<NBLOCK) THEN
エラー処理
ENDIF
ENDDO
:
[1]
DO JJJ=0,1
JJJ=0のとき、↘の通信を行う。
JJJ=1のとき、↙の通信を行う。
DO JJ=JSTA+JJJ*2,JEND,4
DO J=JJ,MIN(JJ+1,JMAX)
DO I=IMIN,IMAX
X(I,J) = ~
ENDDO
ENDDO
ENDDO
:
[2]

```



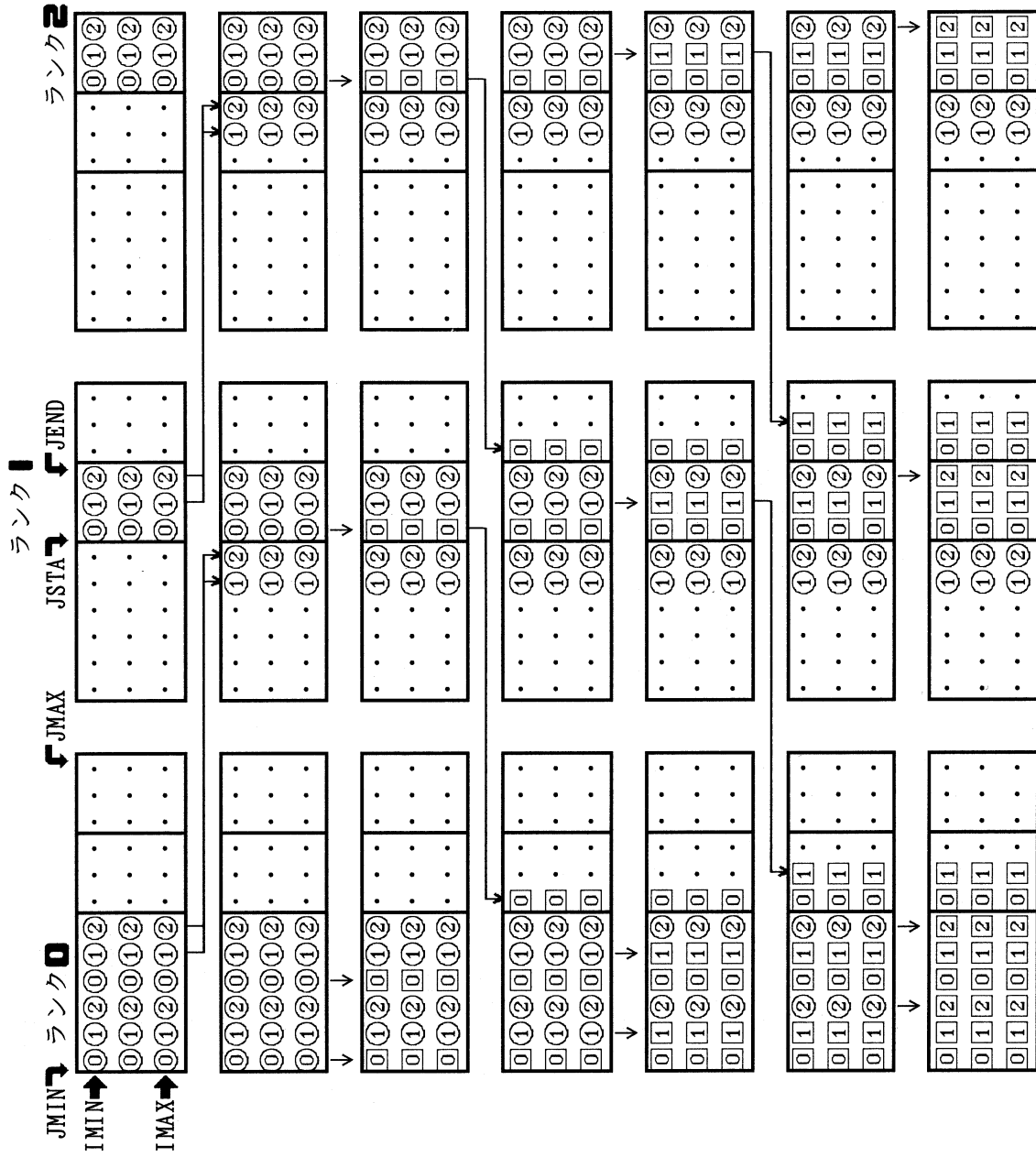
■ ゼブラ (図5-6-10の(3)参照)

- 2次元の13点差分、3次元の25点差分の場合に使用します。
- 各プロセス(ランク**2**を除く)が担当する列数は、3列(①②)の倍数になっています。このため[1]で3を指定し、3列を単位としてブロック分割します。
- 各プロセスが担当する列数は最低3列(①②)必要なので、これを[2]でチェックします。

```

:
NBLOCK = 3
NTOTAL = (JMAX-JMIN)/NBLOCK + 1
DO IRANK=0,NPROCS-1
CALL PARA_RANGE(1,NTOTAL,NPROCS,IRANK,JSTAX,JENDX)
JJSTA(IRANK) = JMIN+(JSTAX-1)*NBLOCK
JJJEND(IRANK) = MIN(JMIN+JENDX*NBLOCK-1,JMAX)
IF (JJJEND(IRANK)-JJSTA(IRANK)+1<NBLOCK) THEN
エラー処理
ENDIF
ENDDO
:
[1]
DO JJ=0,2
JJ=0のとき、↓の通信を行う。
JJ=1のとき、↓の1回目の通信を行う。
JJ=2のとき、↓の2回目の通信を行う。
DO J=JSTA+JJ,JEND,3
DO I=JMIN,IMAX
X(I,J) = ~
ENDDO
ENDDO
ENDDO
:
[2]

```



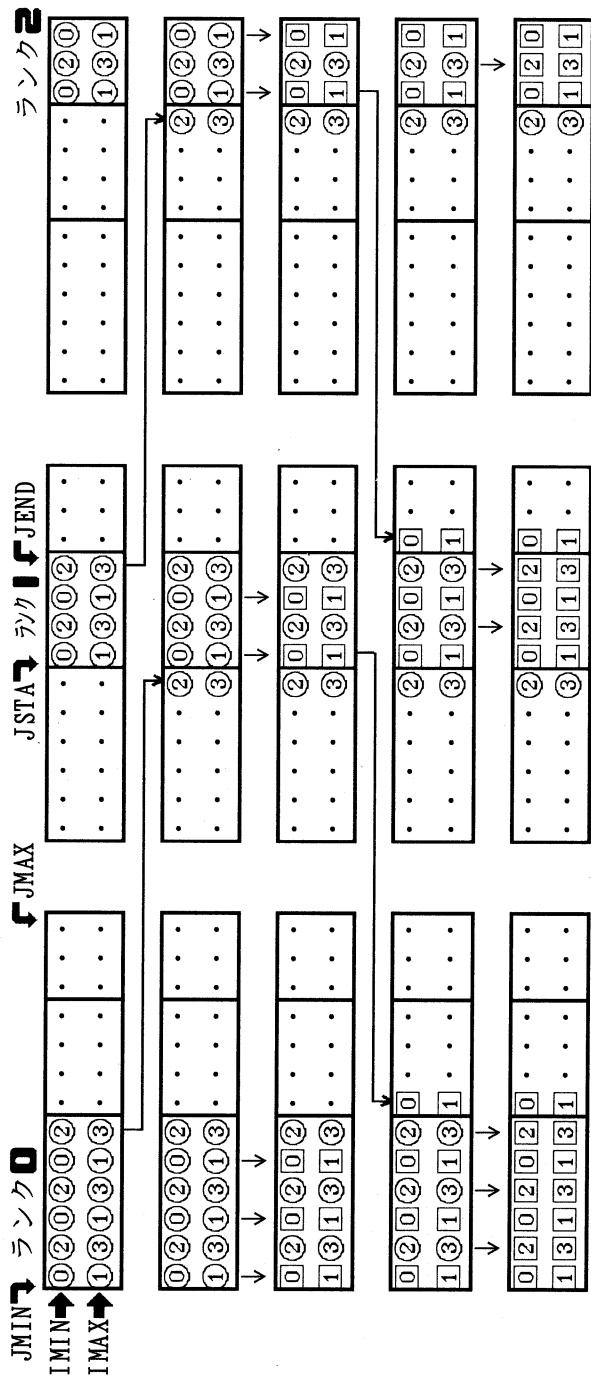
■ マルチカラー (図5-6-10の(4)参照)

- 2次元の5,9点差分、3次元の7,19点差分の場合に使用します。
- 2次元では①～③の(2×2=)4色、3次元では①～⑦の(2×2×2=)8色になります。ただし3次元でJ方向をマルチカラーにしなれば(2×2=)4色です。
- 各プロセス(ランク2)を除く②が担当する列数は、2列(①②)の倍数になっている必要があります。
- このため[1]で2を指定し、2列を単位としてブロック分割します。
- 各プロセスが担当する列数は最低2列(①②)必要なので、これを[2]でチェックします。

```

[1]
NBLOCK = 2
NTOTAL = (JMAX-JMIN)/NBLOCK + 1
DO IRANK=0,NPROCS-1
CALL PARA_RANGE(1,NTOTAL,NPROCS,IRANK,JSTAX,JENDX)
JJSTA(IRANK) = JMIN+(JSTAX-1)*NBLOCK
JJEND(IRANK) = MIN(JMIN+JENDX*NBLOCK-1,JMAX)
IF (JJEND(IRANK)-JJSTA(IRANK)+1<NBLOCK) THEN
  エラー処理
ENDIF
ENDDO
[2]
DO JJ=0,1
  JJ=0のとき、↘の通信を行う。
  JJ=1のとき、↙の通信を行う。
  DO II=0,1
    DO J=JSTA+JJ,JEND,2
      DO I=IMIN+II,IMAX,2
        X(I,J) = ~
      ENDDO
    ENDDO
  ENDDO
ENDDO

```



■ マルチカラー (図5-6-10の(6)参照)

- 2次元の13点差分、3次元の25点差分の場合に使用します。
- 2次元では①～⑧の(3×3=9)色、3次元では①～⑳の(3×3×3=27)色になります。ただし3次元でJ方向をマルチカラーにしなれば(3×3=9)色です。
- 各プロセス(ランク2)を除くが担当する列数は、3列(①③⑤)の倍数になっている必要があります。
- このため[1]で3を指定し、3列を単位としてブロック分割します。
- 各プロセスが担当する列数は最低3列(①③⑤)必要なので、これを[2]でチェックします。

```

NBLOCK = 3
NTOTAL = (JMAX-JMIN)/NBLOCK + 1
DO IRANK=0,NPROCS-1
CALL PARA_RANGE(1,NTOTAL,NPROCS,IRANK,JSTAX,JENDX)
JJSTA(IRANK) = JMIN+(JSTAX-1)*NBLOCK
JJEND(IRANK) = MIN(JMIN+JENDX*NBLOCK-1,JMAX)
IF (JJEND(IRANK)-JJSTA(IRANK)+1<NBLOCK) THEN
  エラー処理
ENDIF
ENDDO

```

[1]

```

DO IRANK=0,2
JJ=0のとき、↘の通信を行う。
JJ=1のとき、↙の1回目の通信を行う。
JJ=2のとき、↙の2回目の通信を行う。
D0 II=0,2
D0 J=JSTA+JJ,JEND,3
D0 I=IMIN+II,IMAX,3
X(I,J) = ~
ENDDO
ENDDO
ENDDO

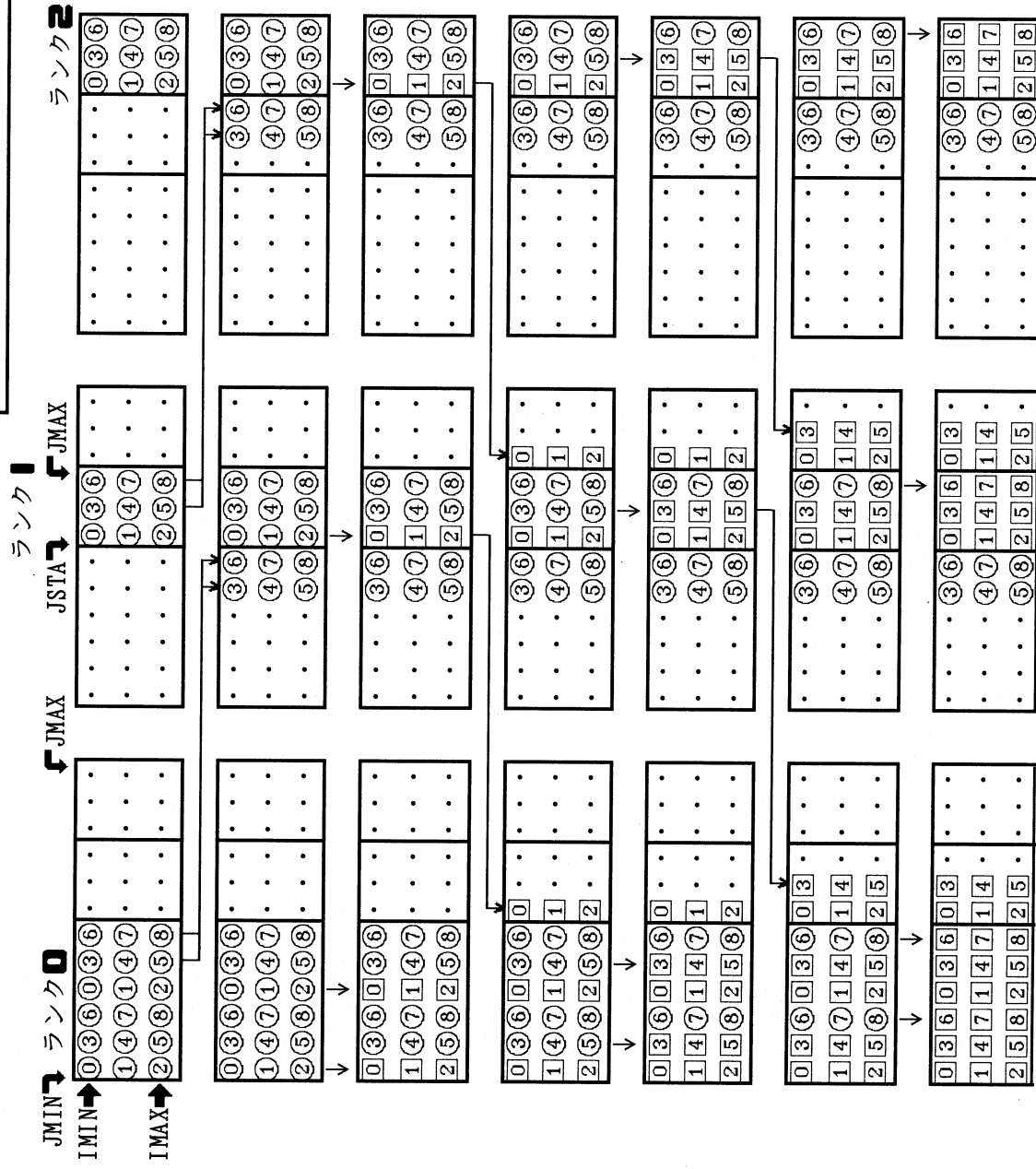
```

[2]

```

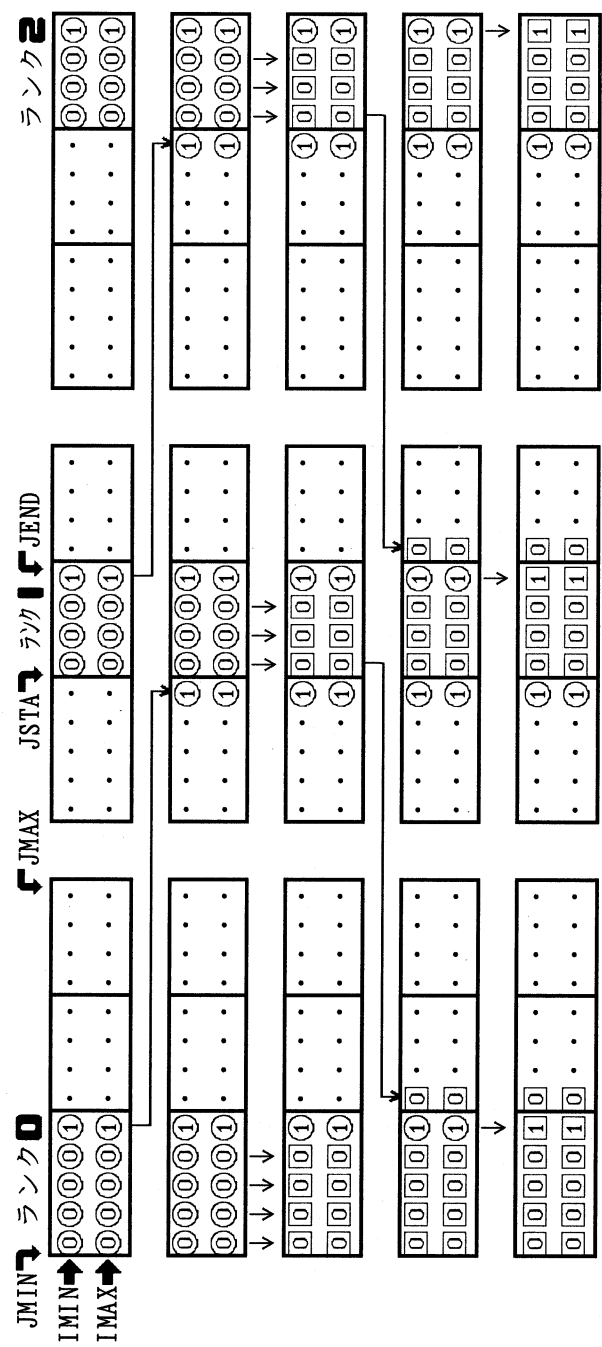
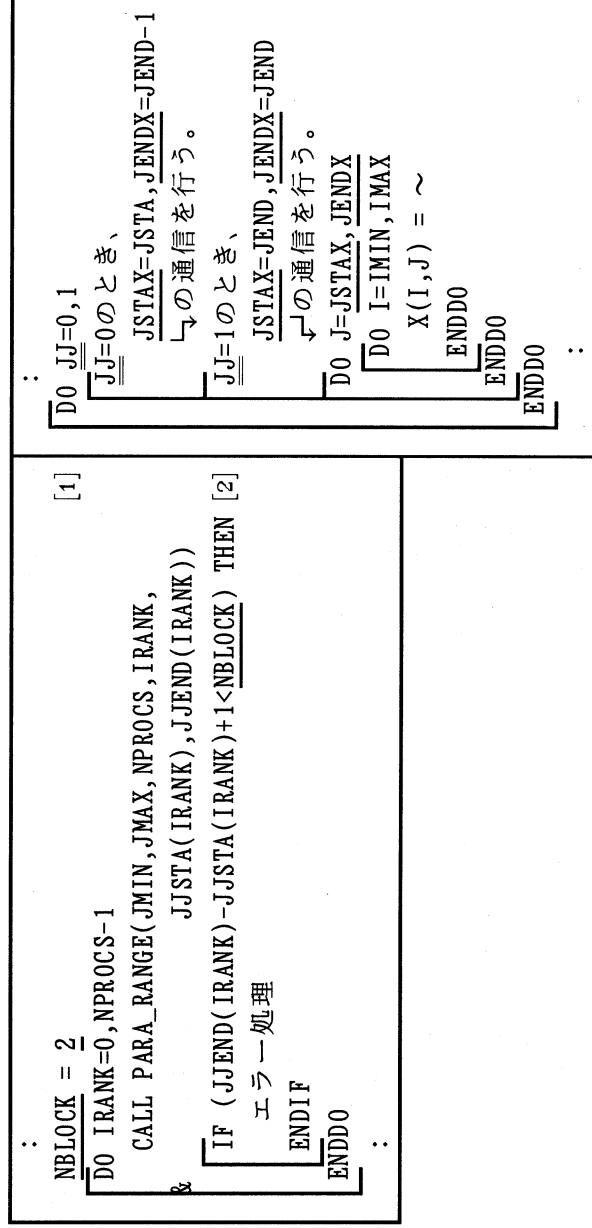
D0 JJ=0,2
JJ=0のとき、↘の通信を行う。
JJ=1のとき、↙の1回目の通信を行う。
JJ=2のとき、↙の2回目の通信を行う。
D0 II=0,2
D0 J=JSTA+JJ,JEND,3
D0 I=IMIN+II,IMAX,3
X(I,J) = ~
ENDDO
ENDDO
ENDDO

```



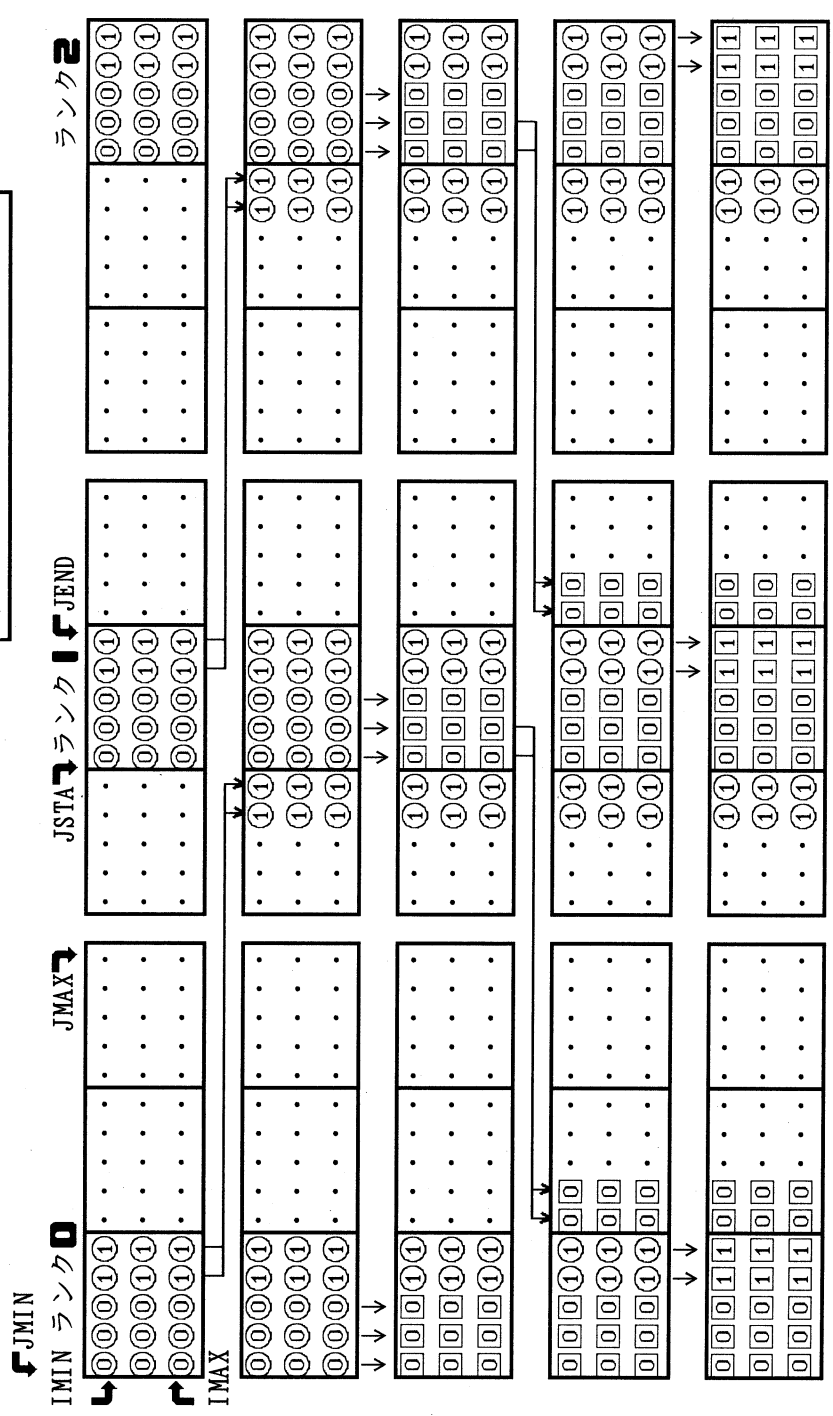
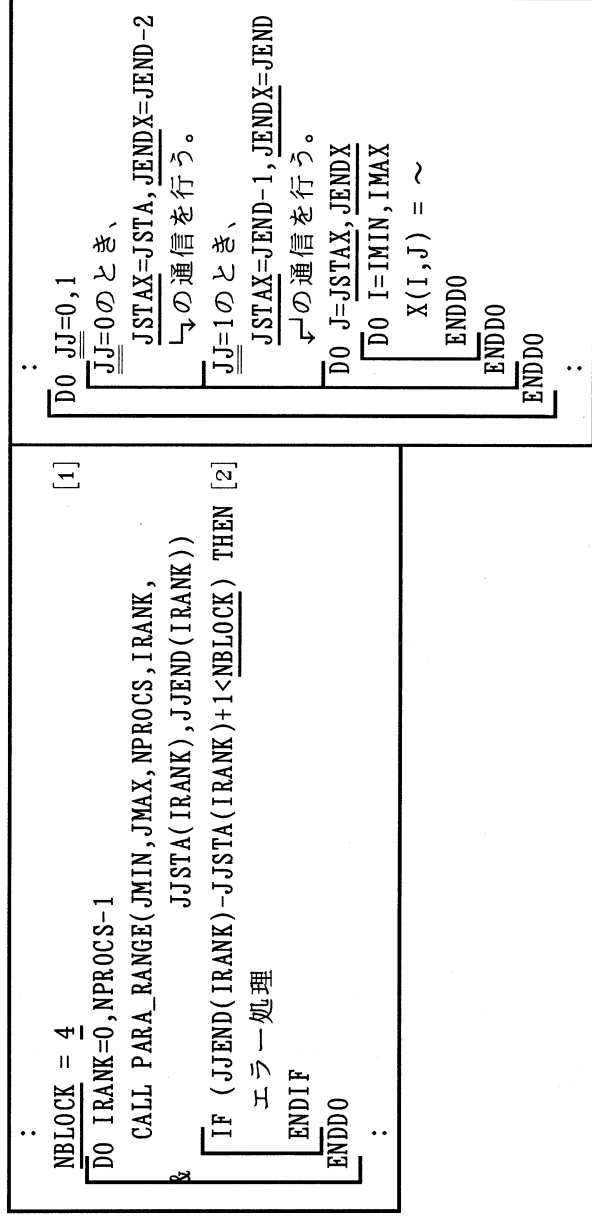
■ パラレル・ブロック・オーダリング (図5-6-10の(7)参照)

- 2次元の5,9点差分、3次元の7,19点差分の場合に使用します。
- 各プロセスが担当する列数は2列以上(ⓐが1列以上、ⓑが1列)必要なので、これを[1][2]でチェックします。



■ パラレル・ブロック・オーダリング (図5-6-10の(8)参照)

- 2次元の13点差分、3次元の25点差分の場合に使用します。
- 各プロセスが担当する列数は4列以上(ⓐが2列以上、ⓑが2列)必要なので、これを[1][2]でチェックします。





## 5-7 モンテカルロ法

モンテカルロ法は、各粒子の挙動がバラバラなので、一般にベクトル計算機よりも並列計算機に向いています。また一般に各粒子間の動作に関係がないことから、計算量に比べて通信量の割合が少なく、並列化の効果が出るが多いようです。本節ではモンテカルロ法モードキの簡単な例を紹介します。

図5-7-1は、粒子が原点から出発して10行程進んだとき、**0~9**のどのレンジに到達したかを調べるプログラムです(例えば図5-7-1の粒子はレンジ**6**に到達しています)。各粒子は1行程で「1単位長さ」進み、進行方向(360度のいずれか)を乱数で決定します。これを100000個の粒子について調べ、**0~9**の各レンジに入った粒子の個数を求めます。

まず単体版のプログラムを図5-7-2に示します。なお、このプログラムでは**0~1**の範囲の**一様乱数**を発生する乱数ルーチンを使用しますが、使用方法はマシン環境によって異なるので具体的には記述しません。

- ①で乱数の種(シード)の初期値を設定します。
- ②のループで100000個の粒子を1つずつ計算します。
- ③で粒子のX, Y座標を「原点」に初期設定します。
- ④のループで粒子を1行程ずつ10行程まで進めます。
- ⑤で乱数を1個発生させ、⑥でそれを使用して粒子の進行方向を決定します。
- ⑦と⑧で粒子を「1単位長さ」進め、X, Y座標を更新します。
- ⑩行程進んだ後の原点との距離を⑨で求め、そのレンジに含まれる個数を⑩で更新します。

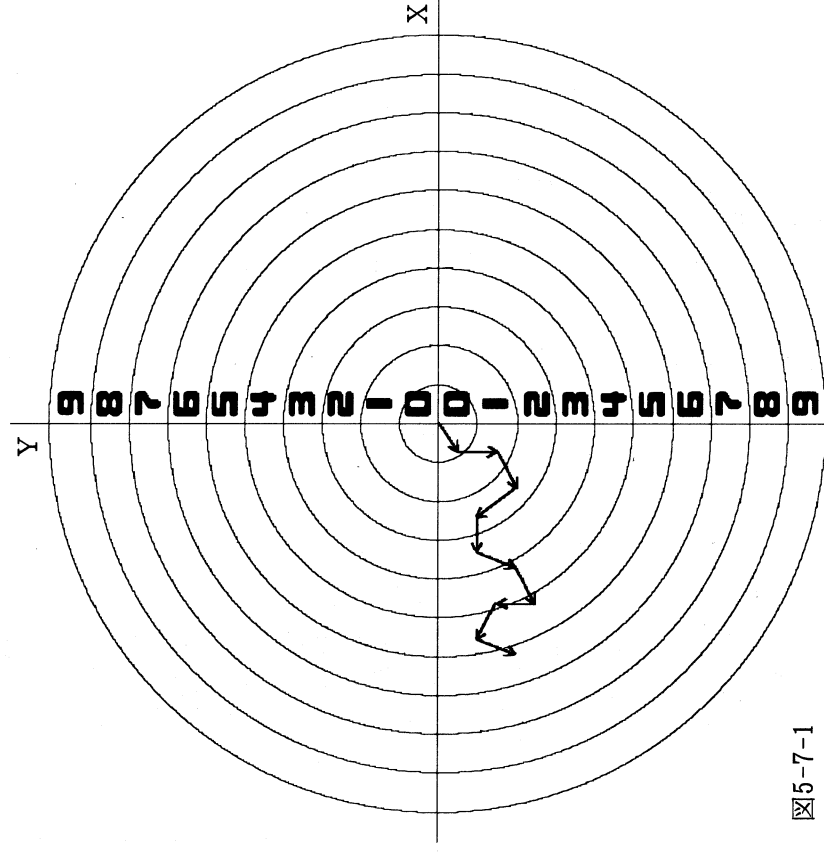


図5-7-1

```

PROGRAM MAIN
PARAMETER (N=100000, PAI=3.1415)
INTEGER ITOTAL(0:9)
DO I = 0, 9
  ITOTAL(I) = 0
ENDDO
乱数の種(シード)の初期値を設定 ①
DO I = 1, N
  X = 0.0; Y = 0.0
  DO ISTEP = 1, 10
    乱数を1個発生させる
    ANGLE = 2.0*PAI*乱数
    X = X + COS(ANGLE)
    Y = Y + SIN(ANGLE)
  ENDDO
  ITEMP = SQRT(X**2 + Y**2) ⑨
  ITOTAL(ITEMP)=ITOTAL(ITEMP)+1⑩
ENDDO
PRINT *, 'TOTAL = ', ITOTAL
END

```

図5-7-2

図5-7-2を並列化したプログラムを図5-7-3に示します。

- 図5-7-3の⑪で、サブルーチン**PARA\_RANGE**(4-5-4節参照)を使用して粒子数(N)をブロック分割し、各プロセスが担当する粒子番号の下限を**ISTA**、上限を**IEND**とします。
- ⑫で乱数の種(シード)の初期値を設定します(詳細は後述します)。
- ⑬で各プロセスは自分の担当する粒子(**ISTA**から**IEND**まで)を並列に計算します。
- ⑭で乱数を1個発生させます(詳細は後述します)。
- これ以後の動作は単体版(図5-7-2)の場合と同じなので、説明は省略します。
- ⑮で、各プロセスが求めた配列**ITOTAL**の値を、**MPI\_REDUCE**を使用して配列の各要素ごとに合計します。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
PARAMETER (N=100000, PAI=3.1415)
INTEGER ITOTAL(0:9), IITOTAL(0:9)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, ~)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, ~)
CALL PARA_RANGE(1, N, NPROCS, MYRANK, ISTA, IEND)
DO I = 0, 9
  IITOTAL(I) = 0
ENDDO
乱数の種(シード)の初期値を設定
DO I = ISTA, IEND
  X = 0.0; Y = 0.0
  DO ISTEP = 1, 10
    乱数を1個発生させる
    ANGLE = 2.0*PAI*乱数
    X = X + COS(ANGLE)
    Y = Y + SIN(ANGLE)
  ENDDO
  ITEMP = SQRT(X**2 + Y**2)
  ITOTAL(ITEMP) = ITOTAL(ITEMP) + 1
ENDDO

```

図5-7-3

↑ 右上へ

左下より

```

CALL MPI_REDUCE(ITOTAL, IITOTAL, 10,
& MPI_INTEGER, MPI_SUM, 0, MPI_COMM_WORLD, IERR)
IF (MYRANK==0) PRINT *, 'TOTAL = ', IITOTAL
CALL MPI_FINALIZE(IERR)
END

```

単体 (1) (2) (3) (4) (5) (6) (7) (8) (9)

図5-7-4 単体版での乱数

ランク0 (1) (2) (3)  
 ランク1 (1) (2) (3) (4) (5) (6)  
 ランク2 (1) (2) (3) (4) (5) (6) (7) (8) (9)

図5-7-5 並列版【方法1】

ランク0 (1) (2) (3) (4) (5) (6) (7) (8) (9)  
 ランク1 (1) (2) (3) (4) (5) (6) (7) (8) (9) (10)  
 ランク2 (1) (2) (3) (4) (5) (6) (7) (8) (9) (10) (11)

図5-7-6 並列版【方法2】

ランク0 (1) (2) (3)  
 ランク1 [1] [2] [3]  
 ランク2 (A) (B) (C)

図5-7-7 並列版【方法3】

並列版で乱数を発生する方法を3種類紹介します(これ以外に、参考文献[70]に記述があります。またネットや文献で調べれば、もっとよい方法があるかもしれません)。説明を簡単にするため、単体版では図5-7-4に示す(1)~(9)の9個の乱数を順番に使用し、並列版は3プロセスがそれぞれ3個の乱数を使用するとします。

●【方法1】(図5-7-5) 各プロセスは単体版と同じ初期値の乱数を発生させ、そのうち下線部の乱数を使用します。例えばランク**2**のプロセスは、ランク**0**と**1**が乱数を3個ずつ使用することが分かっているので、(1)~(6)の乱数をダミーで発生させた後、(7)~(9)の乱数を使用します。図5-7-3では各粒子が10個の乱数を使用するので、**図**の後には(ISTA-1)\*10個の乱数をダミーで発生させた後、**図**で実際に使う乱数を発生させます。この方法は単体版と同じ乱数と同じ順番に使用するので、計算結果も同じになります。しかし図5-7-5で例えばランク**2**のプロセスは、単体版と同じ数の乱数を発生させる必要があるので、乱数を発生する部分に計算時間がかかっているプログラムの場合は、並列化の効果が出ません。また本例のように、各プロセスが使用する乱数の個数が計算前に分かっているプログラムは少ないと思われま。

●【方法2】(図5-7-6) 各プロセスは単体版と同じ初期値の乱数を発生させ、それを各プロセスがサイクルリックに使用します。まず**図**の後でMYRANK個(ランク**2**なら2個)の乱数をダミーで発生させます。**図**ではNPROCS個(本例では3個)の乱数を発生させ、そのうち最初の乱数のみ(図5-7-6の下線の乱数)を使用します。この方法は単体版と同じ乱数を使用しますが、使用する順番が異なるので、計算結果は変わります。また各プロセスは単体版と同じ数の乱数を発生させる必要があるので、乱数を発生する部分に計算時間がかかっているプログラムの場合は、並列化の効果が出ません。

●【方法3】(図5-7-7) 各プロセスはそれぞれ別の初期値の乱数を使用します。**図**で各プロセスが異なる種(シード)の初期値を設定します(初期値を乱数系列の1周期の1/NPROCS個おきに設定する方法、「MYRANK+100」のようにMYRANKの値を使って設定する方法、その時点の時刻を使って設定する方法などがあります)。

この方法は単体版と異なる乱数を使用して計算を行うので、単体版と計算結果が変わります。また1つの並列ジョブ内で異なる初期値の乱数を複数使用するので、乱数の性質が重要となるプログラムの場合は問題があると思われま。一方、例えば3プロセスで実行する場合、発生させる乱数の数は単体版の1/3に減少するので、乱数を発生する部分に計算時間がかかっているプログラムの場合は、並列化の効果が出ま。

個別要素法や分子動力学法のプログラムは並列化に向いていると言われています。通常、これらのプログラムは粒子で反復する一部のD0ループのみがホットスポットになっています。4-3節の『計算パターン1』で述べたように、ホットスポットのD0ループのみを(本当の意味の)並列化し、ループ終了後に、各プロセスは自分が計算した部分を他の全プロセスに送信する方法を取ると、並列化にともなう修正が簡単になり、またパフォーマンスもそこそこ出るようです。

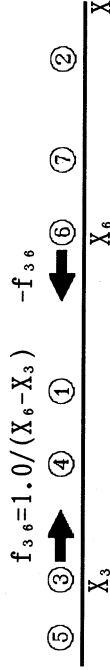
ここでは簡単のため、一次元の個別要素法/分子動力学法モードキの例で並列化の方法を説明します。以下のように一次元の計算領域に粒子が7個存在し、粒子①のX方向の座標を $X_1$ とします。

全ての粒子間には粒子間の距離に反比例した引力が働くとします。例えば粒子③に働く粒子⑥からの引力は  $f_{3,6}=1.0/(X_6-X_3)$  となり、相手の粒子⑥には反対向きの引力 $-f_{3,6}$ が働きます。粒子③に働く全ての粒子からの引力(以後合力と呼びます)を $F_3$ とすると、 $F_3$ は以下ようになります。

$$F_3 = -f_{1,3}-f_{2,3}+f_{3,4}+f_{3,5}+f_{3,6}+f_{3,7}$$

各粒子の合力の計算が終了すると、その値を使用して各粒子の新しい座標を求めます。例えば粒子③に働く合力 $F_3$ と粒子③の現在の座標 $X_3(OLD)$ から、粒子③の新しい座標 $X_3(NEW)$ を次のように求めます。

$$X_3(NEW) = X_3(OLD) + F_3$$



以上の計算を行うプログラムを図5-8-1(1)に示します。各粒子の合力を配列Fに、座標を配列Xに入れます。

- (1)はタイムステップのループです。(2)で合力の配列Fをゼロクリアします。
- (3)と(4)で各粒子の合力を求めます。(3)のIは合力を求めたい粒子の番号①を、(4)のJは相手の粒子の番号③を表します。重複を避けるため(4)のループは『DO J=I+1,N』とします。(5)で粒子①と③の間に働く引力を変数FIJに求め、それを(6)で粒子①の合力に加え、(7)で粒子①の合力から引きます。
- (3)と(4)のループの計算が終了すると、合力の配列Fの中身は図5-8-2(1)のようになります。
- 合力の計算が終了した後、(8)で、各粒子の新しい座標Xを現在の座標Xと合力Fから求めます。

```

:
PARAMETER(N=7)
DIMENSION F(N),X(N)
:
DO I TIME = 1, 100 (1)
DO I = 1, N
F(I) = 0.0 (2)
ENDDO
DO I = 1, N-1 (3)
DO J = I+1, N (4)
FIJ = 1.0/(X(J)-X(I)) (5)
F(I) = F(I) + FIJ (6)
F(J) = F(J) - FIJ (7)
ENDDO
ENDDO
DO I = 1, N
X(I) = X(I) + F(I) (8)
ENDDO
ENDDO
:
    
```

F(1)		+f <sub>1,2</sub> +f <sub>1,3</sub> +f <sub>1,4</sub> +f <sub>1,5</sub> +f <sub>1,6</sub> +f <sub>1,7</sub>
F(2)	-f <sub>1,2</sub>	+f <sub>2,3}+f<sub>2,4}+f<sub>2,5}+f<sub>2,6}+f<sub>2,7</sub></sub></sub></sub></sub>
F(3)	-f <sub>1,3}-f<sub>2,3</sub></sub>	+f <sub>3,4}+f<sub>3,5}+f<sub>3,6}+f<sub>3,7</sub></sub></sub></sub>
F(4)	-f <sub>1,4}-f<sub>2,4}-f<sub>3,4</sub></sub></sub>	+f <sub>4,5}+f<sub>4,6}+f<sub>4,7</sub></sub></sub>
F(5)	-f <sub>1,5}-f<sub>2,5}-f<sub>3,5}-f<sub>4,5</sub></sub></sub></sub>	+f <sub>5,6}+f<sub>5,7</sub></sub>
F(6)	-f <sub>1,6}-f<sub>2,6}-f<sub>3,6}-f<sub>4,6}-f<sub>5,6</sub></sub></sub></sub></sub>	+f <sub>6,7</sub>
F(7)	-f <sub>1,7}-f<sub>2,7}-f<sub>3,7}-f<sub>4,7}-f<sub>5,7}-f<sub>6,7</sub></sub></sub></sub></sub></sub>	

図5-8-2(1)

図5-8-1(1)

図5-8-1(1)のプログラムの並列化について説明します。ここでは説明を簡単にするため、(3)と(4)のループのみ並列化し、(8)のループは並列化しないことにします。

まず分割方法ですが、通常個別要素法や分子動力学法のプログラムでは、各粒子の計算量が異なっていることが多いことと、(3)と(4)のループで実際に計算する部分が図5-8-2(1)の全体ではなく右上三角部分というブロック分割しにくい形状になっていることから、分割方法をサイクリック分割とします。

まず(3)のループでサイクリック分割して並列化したプログラムを図5-8-1(2)に示します。(3)をサイクリック分割して(9)に変更します。(3)と(4)のループが終了した時点では、各プロセスの合力の配列Fの自身は図5-8-2(2)の上図のようになります。(8)のループは並列化しないので、(8)のループに入る前に通信を行って全プロセスが正しい合力Fの値を持つ必要があります。これを(10)で、4-6-6節で紹介した『重ね合せ』を使って行くと、合力Fは図5-8-2(2)の下図のようになります。結果が配列FFに入るので、(8)の配列Fは(11)のようにFFに修正します。

一方(3)と(4)の2重ループの内側でサイクリック分割した場合、図5-8-1(3)のようになります。サイクリック分割に関係する部分を(12)に示します。この場合、(3)と(4)のループが終了した時点で各プロセスの合力の配列Fの自身は図5-8-2(3)のようになります。以後の処理は図5-8-1(2)の場合と同じです。

● 図5-8-1(2)と図5-8-1(3)のどちらがよいかはプロセス間のロードバランス、キャッシュミス、(12)のIF文のオーバーヘッドなどが関係するのでケースバイケースです。

● (3)と(4)の部分の計算量が少ない場合、本例のように重ね合せのような派手な通信を行うと、通信のオーバーヘッドによってパフォーマンスが出ない場合もあります。(3)と(4)部分のチューニング方法については参考文献[16]を参照して下さい)。そのような場合、実際のプログラムでは全粒子間で引力が働くわけではなく(近隣の粒子間のみ働く)ので、重ね合せではなく、必要最小限の要素のみを通信するようにします。ただしプログラムが面倒になります。

● (8)の部分も含めて並列化しないと効果が出ない場合もあります。その場合、タイムステップ内の全ループを検討し、どこで何をどのように通信すれば通信を少なくできるのかを検討する必要があります。

```

:
INCLUDE 'mpif.h'
PARAMETER(N=7)
DIMENSION F(N), X(N), FF(N)
:
DO ITIME = 1, 100
DO I = 1, N
F(I) = 0.0
ENDDO
DO I = 1+MYRANK, N-1, NPROCS (9)
DO J = I+1, N
FIJ = 1.0/(X(J)-X(I))
F(I) = F(I) + FIJ
F(J) = F(J) - FIJ
ENDDO
ENDDO
CALL MPI_ALLREDUCE(F, FF, N, MPI_REAL, (10)
MPI_SUM, MPI_COMM_WORLD, IERR) (10)
DO I = 1, N
X(I) = X(I) + FF(I) (11)
ENDDO
ENDDO
:
&

```

図5-8-1(2)

```

:
INCLUDE 'mpif.h'
PARAMETER(N=7)
DIMENSION F(N), X(N), FF(N)
:
DO ITIME = 1, 100
DO I = 1, N
F(I) = 0.0
ENDDO
IRANK = -1
DO I = 1, N-1
DO J = I+1, N
IRANK = IRANK + 1 (12)
IF (IRANK == NPROCS) IRANK = 0 (12)
IF (MYRANK == IRANK) THEN (12)
FIJ = 1.0/(X(J)-X(I))
F(I) = F(I) + FIJ
F(J) = F(J) - FIJ
ENDIF (12)
ENDDO
ENDDO
CALL MPI_ALLREDUCE(F, FF, N, MPI_REAL,
MPI_SUM, MPI_COMM_WORLD, IERR)
&

```

```

DO I = 1, N
  X(I) = X(I) + FF(I)
ENDDO
ENDDO
  
```

図5-8-1(3)

【ランク0】

F(1)	$+f_{12}+f_{13}+f_{14}+f_{15}+f_{16}+f_{17}$
F(2)	$-f_{12}$
F(3)	$-f_{13}$
F(4)	$-f_{14}$
F(5)	$-f_{15}$
F(6)	$-f_{16}$
F(7)	$-f_{17}$
	$+f_{45}+f_{46}+f_{47}$
	$-f_{45}$
	$-f_{46}$
	$-f_{47}$

【ランク1】

	$+f_{23}+f_{24}+f_{25}+f_{26}+f_{27}$
	$-f_{23}$
	$-f_{24}$
	$-f_{25}$
	$-f_{26}$
	$-f_{27}$
	$+f_{56}+f_{57}$
	$-f_{56}$
	$-f_{57}$

【ランク2】

	$+f_{34}+f_{35}+f_{36}+f_{37}$
	$-f_{34}$
	$-f_{35}$
	$-f_{36}$
	$-f_{37}$
	$+f_{67}$
	$-f_{67}$

図5-8-2(2)

FF(1)	$+f_{12}+f_{13}+f_{14}+f_{15}+f_{16}+f_{17}$
FF(2)	$+f_{23}+f_{24}+f_{25}+f_{26}+f_{27}$
FF(3)	$+f_{34}+f_{35}+f_{36}+f_{37}$
FF(4)	$+f_{45}+f_{46}+f_{47}$
FF(5)	$+f_{56}+f_{57}$
FF(6)	$+f_{67}$
FF(7)	
	$-f_{15}-f_{25}-f_{35}-f_{45}$
	$-f_{16}-f_{26}-f_{36}-f_{46}-f_{56}$
	$-f_{17}-f_{27}-f_{37}-f_{47}-f_{57}-f_{67}$

【ランク0】

F(1)	$+f_{12}$
F(2)	$+f_{23}$
F(3)	$+f_{35}$
F(4)	$+f_{45}$
F(5)	$-f_{35}-f_{46}$
F(6)	$-f_{26}$
F(7)	$-f_{56}$

【ランク1】

	$+f_{13}$
	$+f_{24}$
	$+f_{36}$
	$+f_{46}$
	$+f_{57}$
	$-f_{13}$
	$-f_{24}$
	$-f_{36}-f_{46}$
	$-f_{27}$
	$-f_{57}$

【ランク2】

	$+f_{14}$
	$+f_{25}$
	$+f_{37}$
	$+f_{47}$
	$-f_{14}$
	$-f_{34}$
	$-f_{25}$
	$-f_{17}$
	$-f_{37}-f_{47}$
	$-f_{67}$

図5-8-2(3)

	$+f_{15}$
	$+f_{26}$
	$+f_{35}$
	$+f_{45}$
	$-f_{35}-f_{46}$
	$+f_{56}$
	$-f_{26}$
	$-f_{56}$

(このページは空白です。)

## 第 6 章

# 並列版数値計算ライブラリー

本章では、並列版の数値計算サブルーチンライブラリーについて簡単に紹介します。

## 6-1 並列版数値計算ライブラリー

並列版(分散メモリー用)の数値計算ライブラリーには次の種類があります。(2)と(3)のURLは巻末の参考文献[42]を参照して下さい。

- (1) コンピュータメーカーが提供している数値計算ライブラリー(一般に、そのメーカーのマシン用にチューニングされています)ので高速です)。
- (2) ソフトウェア会社が販売している数値計算ライブラリー(NAG, IMSLなど)。
- (3) フリーの数値計算ライブラリー(PCP, ScaLAPACK, PETSc, GeofEMなど)。

### ■ 単体版の数値計算ライブラリーを使用した並列化

並列版(分散メモリー用)の数値計算ライブラリーは、各プロセスが配列を縮小して持つため、一般に使い方が面倒です。プログラムによっては、並列版でなく、単体版の数値計算ライブラリーを使用して並列化できる場合もあり、一般に後者の方が修正が簡単です。

【例1】行列乗算の数値計算ライブラリーの単体版と並列版が提供されているとします。例えば  $AB=C$  を計算する場合、各プロセスが図6-1-1の部分を計算するようにすれば、単体版のライブラリーを使用して並列化することができます。ただし全プロセスが配列Aの同じ内容を持つ必要があります。

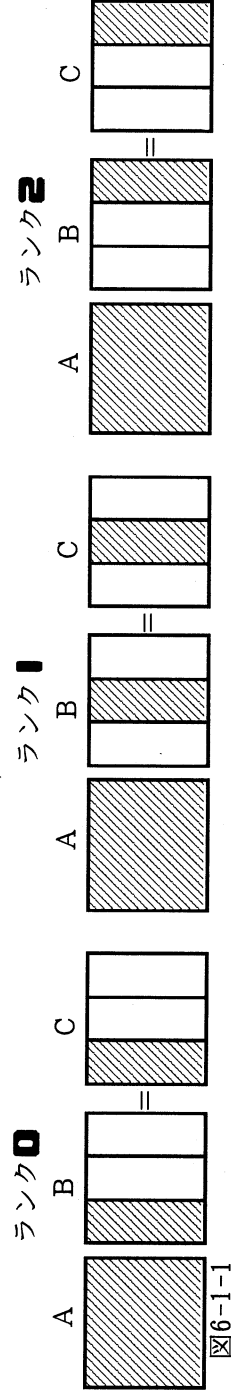


図6-1-1

【例2】例えば6組の非対称密行列の連立一次方程式  $A_i x_i = b_i$  ( $i=1\sim 6$ ) を解くプログラムがあり、それぞれの連立一次方程式は全く独立に解くことができます。連立一次方程式をLU分解で解く数値計算ライブラリーの単体版と並列版が提供されているとします。

並列版のライブラリーでは、図6-1-2(1)のように、各  $A_i x_i = b_i$  ごとに、全プロセスが並列に計算を行います。このとき図の⇨に示す通信が発生します。

一方、図6-1-2(2)のように、6組の連立一次方程式を2組ずつ各プロセスに分配すれば、単体版のライブラリーを使用して並列化することができます。これは4-3節で述べた計算パターン3(上位の部分の並列化)に相当します。

両者を比較すると、図6-1-2(1)では⇨に示す通信が発生するなどの理由で、通常、図6-1-2(2)の方がよいパフォーマンスが得られます。また【例1】と同様の理由で、一般に単体版の方が使用方法が簡単です。

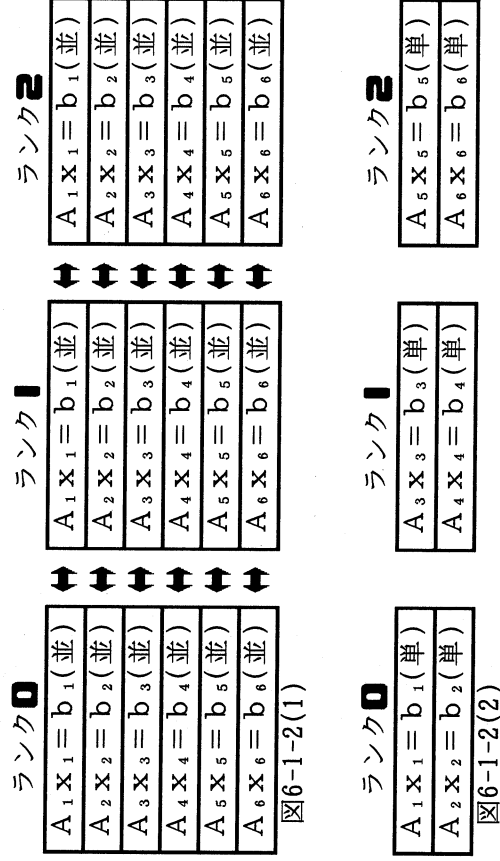


図6-1-2(1)

図6-1-2(2)



## 第7章

### 本当に並列化の必要があるか？

本書もついに最後の章に到達しました。せっかく並列化について勉強してきたところに水をさすようです縮ですが、本章では『本当に並列化の必要があるか？』といういささか過激なタイトルで説明を行います。

## 7-1 単体チューニング

### ■ 単体チューニングと並列化

本書の、しかも7章まで読み進んできた方ならば、何らかの形で並列プログラムを作成してみたいと考えられると思います。そしてその主な目的は経過時間の短縮にあると思われま

す。ところで、経過時間の短縮ということであれば、並列化以外に単体チューニング(参考文献[6]参照)という方法もあります。単体チューニングを行っても全然速くならないプログラムも勿論ありますが、極端なプログラムになると10倍、あるいは100倍速くなる(言いかえると元のプログラムが悪過ぎる)場合があります。

並列化も単体チューニングも、経過時間の短縮という意味では同じことなのですが、根本的に異なる点があります。例えば図7-1-1で、もともと100分かかっているプログラムを単体チューニングし、その結果50分になったとします。この場合、CPU時間が物理的に50分短縮されたわけですから、そのジョブの経過時間が短縮したのと同時に、他のジョブにとってもCPU時間を使用する機会が増えます。つまり単体チューニングは並列計算機のシステム全体の処理効率を向上させます。

一方、単体チューニングをせずに並列化のみを行い、4ノードで経過時間が30分になったとします。考えしてみると、並列化というのは単にCPU時間を各ノードにばらまいているだけであり、CPU時間自体が減少するわけではありません。言いかえると他のジョブのCPU時間を奪って経過時間を短縮させていることになりま

す。さらに、並列化して4ノードで4倍のパフォーマンスを得ることは通常ありえないので、システム全体のCPU時間の合計は単体で実行するよりも多くなります(図7-1-1の例では、単体版で100分だったのが並列版ではシステム全体で $30 \times 4 = 120$ 分)。このように、並列化というのは並列計算機のシステム全体の処理効率を低下させることになりま

経過時間

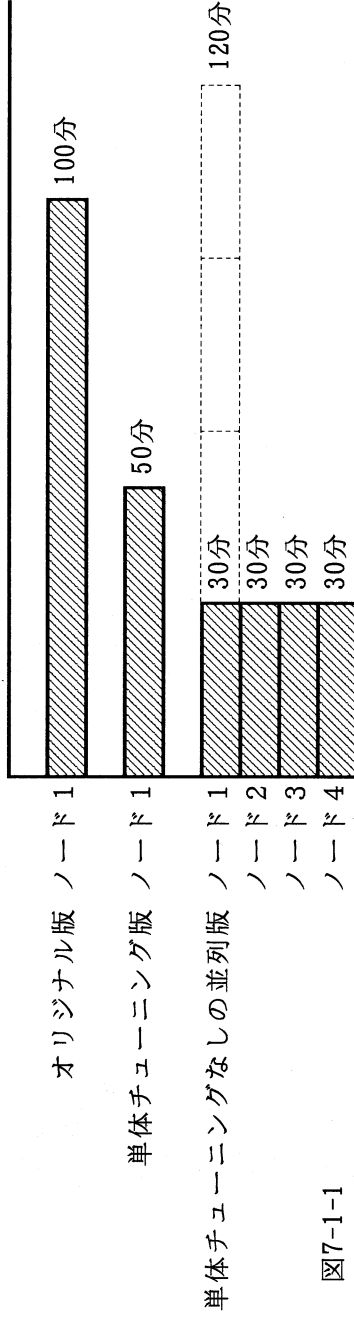


図7-1-1

勿論、単体チューニングだけではパフォーマンス向上に限界があり、それ以上のパフォーマンスを得る必要があるのであれば、やはり並列化するしかありません。いずれにしても、並列化ばかりに目を奪われず、まず単体チューニングを行ってみて、その後で、さらに経過時間を短縮させる必要があれば並列化を行うのがよいのではないかと思います。

チューニングの効果は、特に多くのノードを使用して並列計算する場合、顕著に現れます。図7-1-2に示すように、ある単体プログラムをチューニングして速度が2倍になったとします。このプログラムは100%の部分が並列化でき、通信がほぼゼロだとすると、オリジナル版を20台で並列計算したときの経過時間と、チューニング版を10台で並列計算したときの経過時間が同じになり、チューニング版では10台分のノードが節約できることになりま

経過時間

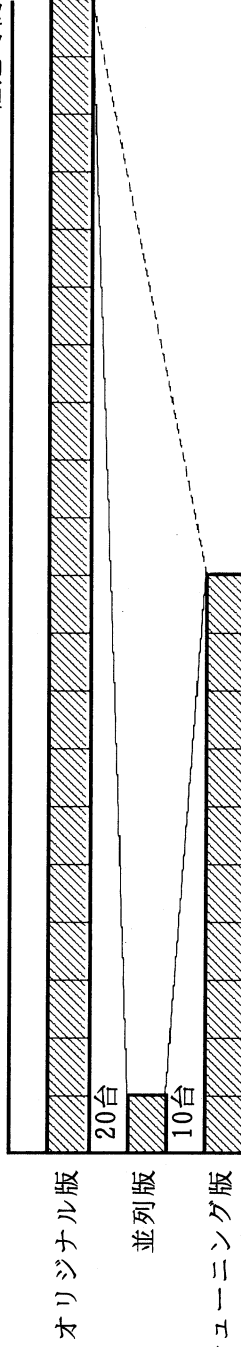


図7-1-2

## 7-2 並列処理と分散処理

4-3節で、プログラムには末端の部分での並列化と上位の部分の並列化があり、なるべく上位の部分で並列化した方がよいという説明をしました。ところで上位の部分の並列化をもっと上に進めていくとどうなるのでしょうか？

### ■ ジョブ単位の並列化

一人のユーザーが計算機を利用して仕事を進めていく場合、大きく次の2つの形態が考えられます。

#### (1) ジョブ間に並列性がない場合

この形態では、データ1を入力としてジョブを実行し、実行結果1が得られます。次に、その結果に基づいてデータ2を作成します。そしてデータ2を入力として2本目のジョブを実行します。以後はこの繰り返しになります。この形態の場合、ジョブ間に並列性はなく、ジョブを4本同時に実行することはできません。



図7-2-1(1)

#### (2) ジョブ間に並列性がある場合

この形態では、データ1~4種類のデータ(例えば少しずつ条件の異なるデータ)を最初に作成し、その後でジョブを実行します。この場合はジョブ間に並列性があるので、4本のジョブを同時に実行することができます。



図7-2-1(2)

図7-2-1(1)の形態の場合、結果4を得るまでの時間を短縮しようとするれば、やはりプログラム自体を並列化する必要があります。

一方図7-2-1(2)の形態の場合、ジョブ間に並列性があります。並列計算機には複数のノードがあるため、最終結果4を得るまでの時間を短縮したければ、プログラム自体を並列化しなくとも、4つのジョブを並列計算機の複数のノードで同時に実行すればよいのです。図7-2-1(2)の形態が、まさにさきほど述べた、1つのプログラムの上位の部分の並列化のさらに上位の並列化、すなわちジョブ単位の並列化であると言えます。通常これを分散処理と呼びます。

ジョブ単位の並列化では、当然ながらプログラムを並列化する必要はなく、また、4ノードで3倍しかパフォーマンスが向上しないといった、並列化にともなうシステム全体の効率の低下もありません。

ジョブ単位の並列化をさらに上に進めていくとどうなるでしょうか？ そうです。最後は人間単位の並列化になります。一般に並列計算機は一人のユーザーだけでなく、多くの人が使用します。このとき、Aさん、Bさん、Cさんのジョブ間には(通常)因果関係はなく、ジョブ間に並列性があります。人間単位の並列化もジョブ単位の並列化の一種なので、以後は両者を合わせてジョブ単位の並列化と呼びます。

## ■ 負荷分散用ソフトウェア

以上のように、大勢の人が並列計算機を使用している環境、または単一ユーザーのジョブ間に並列性がある環境では、同時に実行可能な多くのジョブがあります。一方、並列計算機は多くのノードから構成されており、並列ジョブを実行するマシンであると同時に、多くの単体ジョブを処理する分散処理マシンであるということができます。

ユーザーと並列計算機の間には、一般に負荷分散用ソフトウェアが介在します。その一例を図7-2-2に示します。この例では各ノードが単体ジョブ用と並列ジョブ用に分かれています。まずユーザーは単体ジョブや並列ジョブを負荷分散用ソフトウェアに対して投入します。負荷分散用ソフトウェアはジョブをいったん待ち行列(単体ジョブ用または並列ジョブ用)に入れ、そして待ち行列の中で優先順位の高いジョブを、その時点で最も負荷の軽いノードに割り当てます。

これによって、並列計算機のシステム全体の処理能力を最大にすることができ、また、どのノードが空いているかをユーザーがいちいち調べる手間も省けます。

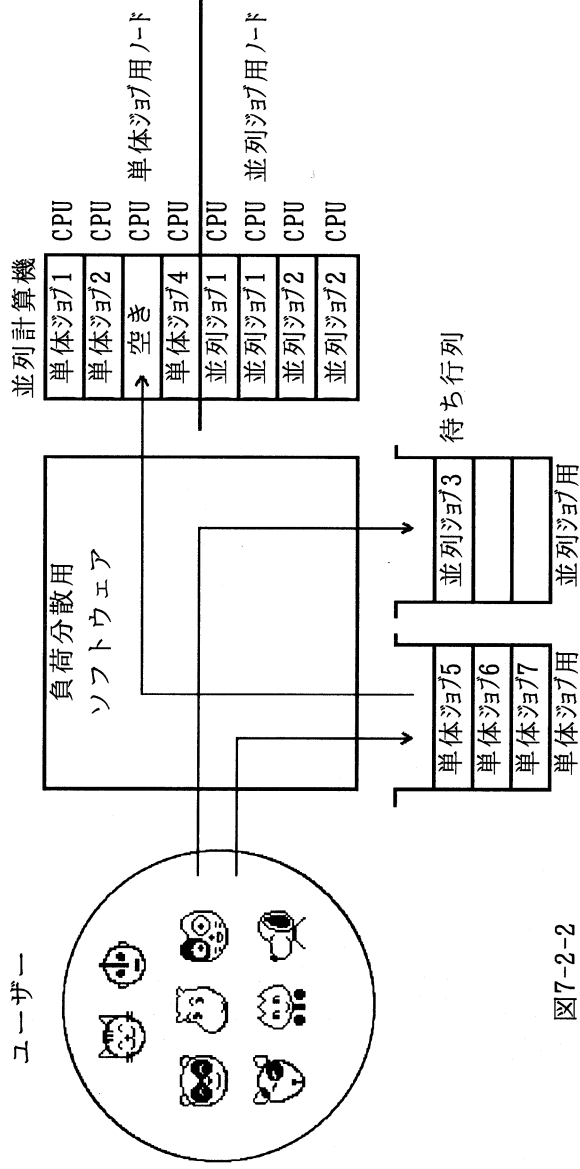


図7-2-2

## ■ 並列計算機と並列ジョブ

いつの間にか、並列化から分散処理に話がそれてしまいました。それでは一体どのような場合にプログラムの並列化を行うのでしょうか？ それは主に以下の2つのケースになると思います。

- 何度も使用し、計算時間が非常に長いので、とにかく巨大なメモリが欲しい場合。
- 計算領域のメモッシュを多くしたいので、とにかく短縮したい場合。

次に、システム運用の観点から見ただけの場合、並列ジョブと分散処理をどのように共存させたらよいのでしょうか？ 並列ジョブは最も遅いノードに足をひっぱられるため、高いパフォーマンスを得るためには、他のジョブが流れていない占有環境で実行するのは難しくなります。従って、多くのユーザーが多くのジョブを実行している環境では並列ジョブを実行するのは難しくなります。このような場合、並列ジョブの実行は夜間と休日のみにするなど時間帯によって制限する方法や、ノード数がたくさんあるのであれば、前述のように並列計算機の各ノードを並列ジョブ用と単体ジョブ用に分ける方法などがあります。

## 7-3 並列化とワークロード

### 第7章 本当に並列化を行う必要があるか？

プログラムを並列化するのにはどの程度ワークロードがかかるのかという質問がよく出ます。本節ではこれに対する回答および、並列化に関するその他のトピックについて説明します。

● 4-3節で説明したように、まず、全体的にどのようなパターンで並列化を行うかを決定する必要があるありますが、そのためにはプログラムの全体構造をある程度把握する必要があります。私の場合は他の人が作成した中身が全くわからないプログラムを並列化するので、実は並列化作業の大部分のワークロードはプログラムの全体構造の把握に費やされています。逆に、プログラムの中身を知っていれば、この時間はそれほどかからないでしょう。

● ベクトル計算機の場合も、ベクトル化に向いているプログラムとそうでないプログラムがありますが、並列化の場合も同様です。並列化に向いていないプログラムを無理に並列化しても意味はありません。

● 並列化できるプログラムの場合、修正箇所が少なく、10分程度で完成するプログラムもあれば、修正箇所が多く、修正作業に時間がかかるプログラムもあります。一般に、4-3節の「計算パターン1」と「計算パターン3」では簡単に並列化を行うことができます。

● 並列化を行うと通信が必要になります。どの変数(配列)をプログラム内のどこで誰に通信すれば、並列プログラムとして正しく動作し、しかも通信量を最小にできるのかを判断する必要がありますが、この判断も、プログラムの中身がある程度知っていない場合はワークロードがかかかります。

● ホットスポットの部分をそのまま並列版の数値計算ライブラリーに置き換えることができれば、並列化のワークロードがあまりかからず効率的で、しかも(一般に)高いパフォーマンスを得ることができます。

● 過去の経験では、並列化の効果がでるのは次のようなプログラムです。

- 差分法・・・ソルバーがSOR法の場合はレッドブランク(マルチカラー)SOR法に修正します。ADI法もパイプライン法またはツイースト分割法で並列化が可能です。ソルバーにICCG法を使用している場合はパイプライン法によって並列化できませんが、メッシュ数が少ないとあまり効果は出ません。

- モンテカルロ法、分子動力学法、個別要素法などの粒子を扱う計算は並列化の効果が出ます。

- 境界要素法・・・ホットスポットが密行列連立一次方程式になるため、並列化の効果が出ます。

逆に、並列化してもあまり効果が出ない解法は、

- 有限要素法(陰解法)・・・疎行列の連立一次方程式をスカイライン法で解く場合、並列化は可能ですが、計算量が少ないために通信の割合が多くなってしまい、あまり効果がでません。また、ICCG法の場合は、不規則疎行列なのでパイプライン法が使えません。一般に、有限要素法の場合は、連立一次方程式の部分を並列化するのでなく、もっと上位のレベルで、領域分割法と呼ばれる並列化を行うのがよいと言われています(参考文献[17])。

以上、長々と説明をしてみました。本書が分散メモリ型並列計算機でのプログラム並列化の一助になれば幸いです。一人でも多くの方がプログラムの並列化に成功することを願ってやみません。お疲れさまでした。

(このページは空白です。)

## 付録

### MPI サブルーチン一覧

本付録では、MPIで提供されている主なサブルーチンの使用方法について説明します。また付録内の「機能」の「関連する節」の所に、各サブルーチンが出てくる本書内の節を示しました。

なお、本書は正式マニュアルではありませんので、正式な使用方法は参考文献[15][27]を参照して下さい。

## 機能 (関連する節: 3 - 2 節)

MPI環境の初期化処理を行います。ユーザーは、MPI\_INITが行う処理内容について意識する必要はありません。MPI\_INITは、全てのMPIルーチンの一番最初に1回だけ、必ずコールする必要があります(必須)。MPI\_INITをコールする前にユーザープログラムが存在しても構いません。

## 使用法

```
CALL MPI_INIT(ierr)
```

- `ierror` : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## サンプルプログラム

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  CALL MPI_INIT(ierr)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, ierr)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, ierr)
  PRINT *, 'NPROCS = ', NPROCS, 'MYRANK = ', MYRANK
  CALL MPI_FINALIZE(ierr)
END
```

MPIが使用するインクルードファイルを指定します。  
MPI環境の初期化処理を行います。  
プロセス数NPROCSを取得します。  
自分のランクMYRANKを取得します。  
NPROCSとMYRANKの値を書き出します。  
MPI環境の終了処理を行います。

```
[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
NPROCS = 3 MYRANK = 0
NPROCS = 3 MYRANK = 1
NPROCS = 3 MYRANK = 2
```



### 機能 (関連する節: 3 - 2 節)

MPI環境の終了処理を行います。ユーザーは、MPI\_FINALIZEが行う処理内容について意識する必要はありません。MPI\_FINALIZEは、全てのMPIルーチンの一番最後に1回だけ、必ずコールする必要があります(必須)。MPI\_FINALIZEをコールした後ユーザープログラムが存在しても構いません。

プログラムが終了する全ての箇所(STOP文またはEND文)で、必ず1度だけMPI\_FINALIZEをコールしてから終了するようにして下さい。例えば図1では、MPI\_FINALIZEを図2のように挿入します。

なお、あるプロセスがMPI\_FINALIZEをコールせずにSTOP文やEND文で終了した場合、マシン環境によっては誤動作する(実行中の他のプロセスを強制的に終了させてしまう、あるいは実行するたびに計算結果が異なる再現性のないエラーとなる)ことがありますので、必ずコールしてから終了するようにして下さい。

```

:
IF (xxx) STOP
:
CALL MPI_BCAST(~) 最後のMPIルーチン
:
END

```

図1 ✕

```

:
IF (xxx) THEN
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
:
CALL MPI_BCAST(~) 最後のMPIルーチン
CALL MPI_FINALIZE(IERR)
:
END

```

図2 ○

### 使用法

```
CALL MPI_FINALIZE(ierrord)
```

- `ierrord`: 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

### サンプルプログラム

```

PROGRAM MAIN
INCLUDE 'mpif.h'
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR) MPIが使用するインクルードファイルを指定します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR) MPI環境の初期化処理を行います。
PRINT *, 'NPROCS = ', NPROCS, 'MYRANK = ', MYRANK  NPROCSとMYRANKの値を書き出します。
CALL MPI_FINALIZE(IERR) MPI環境の終了処理を行います。
END

```

```

[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
NPROCS = 3 MYRANK = 0
NPROCS = 3 MYRANK = 1
NPROCS = 3 MYRANK = 2

```

# MPI\_COMM\_SIZE

## 機能 (関連する節: 3-2 節)

コミュニケータ `comm` で指定したグループに含まれるプロセスの数が `size` に戻ります。

## 使用法

```
CALL MPI_COMM_SIZE(comm, size, ierror)
```

- `comm` : 整数。コミュニケータを指定します。
- `size` : 整数。`comm` で指定したグループ内に含まれるプロセスの数が戻ります。
- `ierror` : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## サンプルプログラム

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  PRINT *, 'NPROCS = ', NPROCS, 'MYRANK = ', MYRANK
  CALL MPI_FINALIZE(IERR)
END
```

MPIが使用するインクルードファイルを指定します。

MPI環境の初期化処理を行います。

CALL MPI\_COMM\_SIZE(MPI\_COMM\_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。

CALL MPI\_COMM\_RANK(MPI\_COMM\_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。

PRINT \*, 'NPROCS = ', NPROCS, 'MYRANK = ', MYRANK NPROCSとMYRANKの値を書き出します。

CALL MPI\_FINALIZE(IERR) MPI環境の終了処理を行います。

```
[aoyama@node01]~/aoyama: mpiexec -np 3 a.out
NPROCS = 3 MYRANK = 0
NPROCS = 3 MYRANK = 1
NPROCS = 3 MYRANK = 2
```

機能 (関連する節: 3-2 節)

コミュニケータ `comm` で指定したグループ内での自分(コールしたプロセス)のランクが `rank` に戻ります。

使用法

```
CALL MPI_COMM_RANK(comm, rank, ierror)
```

- `comm` : 整数。コミュニケータを指定します。
- `rank` : 整数。`comm` で指定したグループ内での自分(コールしたプロセス)のランクが戻ります。
- `ierror` : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

サンプルプログラム

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  PRINT *, 'NPROCS = ', NPROCS, 'MYRANK = ', MYRANK
  CALL MPI_FINALIZE(IERR)
END
```

MPI が使用するインクルードファイルを指定します。  
MPI 環境の初期化処理を行います。  
プロセス数 `NPROCS` を取得します。  
自分のランク `MYRANK` を取得します。  
`NPROCS` と `MYRANK` の値を書き出します。  
MPI 環境の終了処理を行います。

```
[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
NPROCS = 3 MYRANK = 0
NPROCS = 3 MYRANK = 1
NPROCS = 3 MYRANK = 2
```

# 環境管理サブルーチン MPI\_ABORT

## 機能 (関連する節: 4-4-1節, 4-8-1節)

いずれかのプロセスが「MPI\_ABORT」をコールすると、コミュニケーションに所属する全てのプロセスを強制的に異常終了させます。

## 使用法

```
CALL MPI_ABORT(comm, errorcode, ierror)
```

- **comm** : 整数。異常終了させたいプロセスが所属するコミュニケータを指定します。
- **errorcode** : 整数。ここに指定した値がOSに戻ります(FortranでSTOP文とともに指定する数字に相当します)。この値が異常終了メッセージ内でのように表示されるか(あるいは表示されないか)はマシン環境によって異なりますので、マニュアルを参照して下さい。
- **ierror** : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## サンプルプログラム

```
PROGRAM MAIN
INCLUDE 'mpif.h'
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
:
IF (エラーを発見したら)
PRINT *, '===ERROR=== MYRANK=', MYRANK
CALL MPI_ABORT(MPI_COMM_WORLD, 999, IERR)
CALL MPI_FINALIZE(IERR)
STOP
ENDIF
:
CALL MPI_FINALIZE(IERR)
END
```

MPIが使用するインクルードファイルを指定します。  
MPI環境の初期化処理を行います。  
プロセス数NPROCSを取得します。  
自分のランクMYRANKを取得します。

いずれかのプロセスがエラーを発見したら、  
エラーメッセージとランク番号を表示し、  
全プロセスを強制的に異常終了させます。

MPI環境の終了処理を行います。

## 環境管理ファンクション MPI\_WTIME

## ■ 機能 (関連する節: 4-8-3 節)

『経過時間を測定したい部分』の前後で本ファンクションを実行すると、得られた値の差が、その部分の経過時間(単位は秒)となります。なお、得られた経過時間は、本ファンクションを実行したプロセスでの経過時間であり、他のプロセスとは関係ありません。

## ■ 使用法

```
elp = MPI_WTIME()
```

● elp : 倍精度実数。ある過去の時点からの経過時間(単位は秒)が戻ります。

## ■ サンプルプログラム

以下の例で、『計算部分』の経過時間の測定を、②で開始し④で終了します。⑤に示すように、得られた時刻の差が経過時間(単位は秒)となります。本ファンクションで得られる値は倍精度なので、⑥を宣言して下さい。

なお、①(通常は必須)と③(任意)を指定する意味は 4-8-3 節を参照して下さい。

```
PROGRAM MAIN
```

```
INCLUDE 'mpif.h'
```

```
REAL*8 ELP1,ELP2
```

```
CALL MPI_INIT(IERR)
```

```
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR) プロセス数NPROCSを取得します。。
```

```
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。
```

入力データの読み込みなど

```
CALL MPI_BARRIER(MPI_COMM_WORLD,IERR) ① プロセス間の同期をとります。
```

```
ELP1 = MPI_WTIME() ② MPI_WTIME()を実行して測定を開始します。
```

計算部分

```
CALL MPI_BARRIER(MPI_COMM_WORLD,IERR) ③ (任意)プロセス間の同期をとります。
```

```
ELP2 = MPI_WTIME() ④ MPI_WTIME()を実行して測定を終了します。
```

```
PRINT *,'ELAPSE = ',ELP2-ELP1 ⑤ 経過時間を秒で表します。
```

```
CALL MPI_FINALIZE(IERR) MPI環境の終了処理を行います。
```

```
END
```

## 集団通信サブルーチン MPI\_BCAST

■ 機能 (関連する節: 3-3-4節, 4-6-5節, 5-3節)

コミュニケータ(comm)内の、1つの送信元プロセス(root)の送信バッファ(buffer)から、その他全てのプロセスの受信バッファ(buffer)にメッセージを送信します。

### ■ 使用法

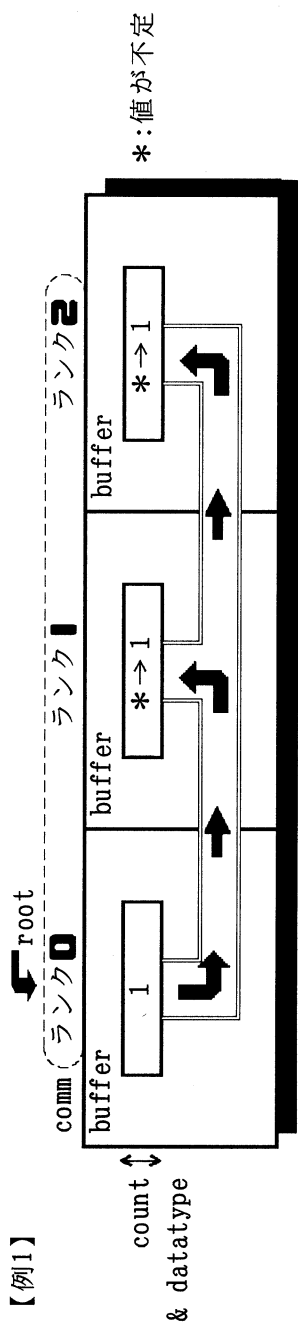
```
CALL MPI_BCAST(buffer, count, datatype, root, comm, ierror)
```

- **buffer** : 送信元プロセスでは、送信バッファの先頭アドレスを指定します。  
宛先プロセスでは、受信バッファの先頭アドレスを指定します。
- **count** : 整数。メッセージの要素数を指定します。
- **datatype** : 整数。メッセージのデータ型を指定します。
- **root** : 整数。送信元プロセスのcomm内でのランクを指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- **comm** : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケータを指定します。
- **ierror** : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

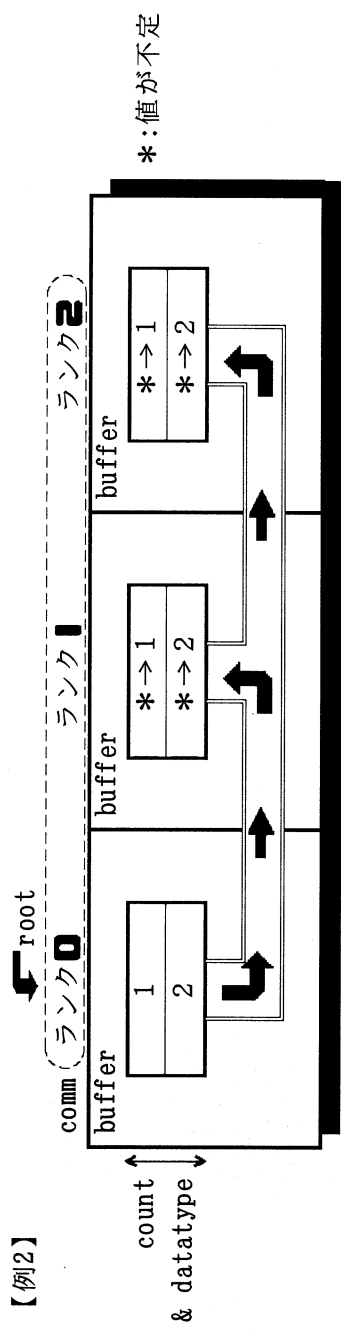
### ■ 注意

- 送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでcount, datatypeは異なっていても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要がある場合があります。

【例1】



【例2】



■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。
IF (MYRANK == 0) THEN
  IBUF = 1
ENDIF
CALL MPI_BCAST(IBUF, 1, MPI_INTEGER, 0, MPI_COMM_WORLD, IERR) 通信を行います。
PRINT *, 'MYRANK = ', MYRANK, 'IBUF = ', IBUF 通信後のバッファの内容を書き出します。
CALL MPI_FINALIZE(IERR)
END
  
```

ランク0のプロセスは送信バッファにデータをセットします。

```

[aoyama@node01]/u/aoyama: mpiexec -np 3 a.out
MYRANK = 0 IBUF = 1
MYRANK = 1 IBUF = 1
MYRANK = 2 IBUF = 1
  
```

# 集団通信サブルーチン MPI\_SCATTER

## 機能 (関連する節: 3 - 3 節)

- コミュニケータ(comm)内の、1つの送信元プロセス(root)の送信バッファ(sendbuf)から、全プロセスの受信バッファ(recvbuf)にメッセージを送信します。
- 各宛先プロセスへの送信メッセージの長さは一定で、送信バッファ(sendbuf)の先頭から宛先プロセスのランクが小さい順に送信されます。
- 本サブルーチンは、MPI\_GATHERと逆の通信になります。

## 使用法

```
CALL MPI_SCATTER(sendbuf, sendcount, sendtype,  
                recvbuf, recvcnt, root, comm, ierr)
```

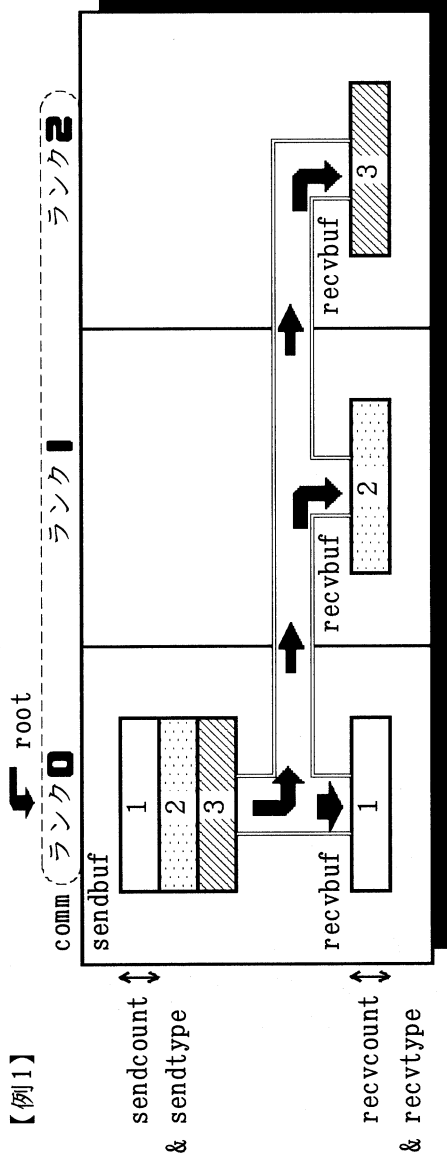
- sendbuf : 送信バッファの先頭アドレスを指定します。rootプロセスのみで意味を持ちます。
- sendcount : 整数。1つのプロセスに送信する送信メッセージの要素数を指定します。  
rootプロセスのみで意味を持ちます。
- sendtype : 整数。送信メッセージのデータ型を指定します。rootプロセスのみで意味を持ちます。
- recvbuf : 受信バッファの先頭アドレスを指定します。  
送信バッファと受信バッファの実際に使用する部分は、メモリー上で重なってははいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1節を参照)。
- recvcnt : 整数。受信メッセージの要素数を指定します。
- recvtpe : 整数。受信メッセージのデータ型を指定します。
- root : 整数。送信元プロセスのcomm内でのランクを指定します。  
comm内の全プロセスが同じ値を指定する必要があるあります。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケータを指定します。
- ierr : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

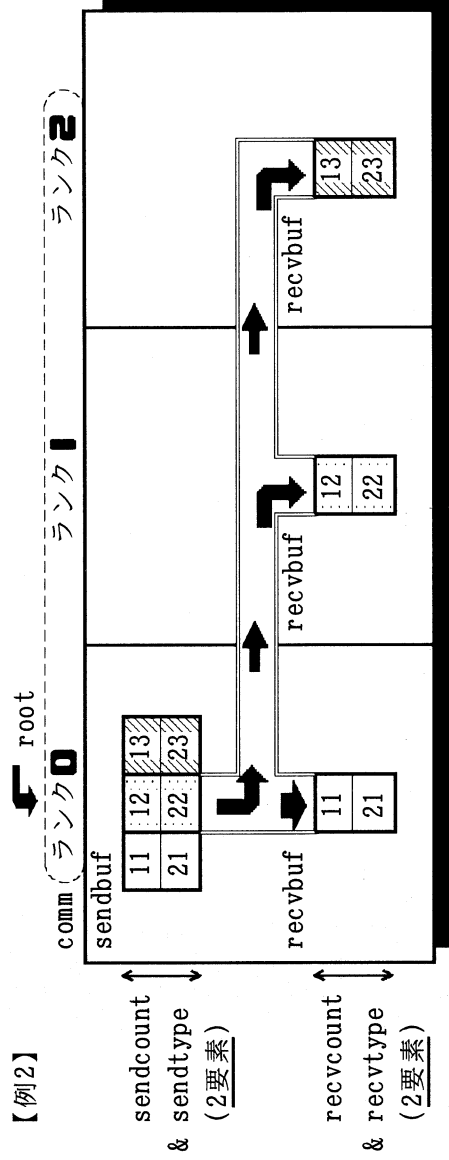
- 各プロセスへの送信メッセージ間、および、送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでsendcount, sendtype, recvcnt, recvtpeは異なっていても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- 送信元プロセス(root)以外のプロセスでは、sendbuf, sendcount, sendtypeは無視されますので、適当な変数を指定してもかまいません。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。



【例1】



【例2】



■ サンプルプログラム 【例1】のプログラム例です。

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISEND(3)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。
  IF (MYRANK == 0) THEN
    DO I = 1, NPROCS
      ISEND(I) = I
    ENDDO
  ENDIF
  CALL MPI_SCATTER(ISEND,1,MPI_INTEGER,
    & IRECV,1,MPI_INTEGER,0,MPI_COMM_WORLD,IERR) 通信を行います。
  PRINT *, 'MYRANK = ',MYRANK, 'IRECV = ',IRECV 通信後の受信バッファの内容を書き出します。
  CALL MPI_FINALIZE(IERR)
  END
```

MPIが使用するインクルードファイルを指定します。  
送信バッファを指定します。

MPI環境の初期化処理を行います。  
プロセス数NPROCSを取得します。  
自分のランクMYRANKを取得します。

```
IF (MYRANK == 0) THEN
  DO I = 1, NPROCS
    ISEND(I) = I
  ENDDO
ENDIF
```

ランク0のプロセスは送信バッファにデータを  
セットします。

通信を行います。

通信後の受信バッファの内容を書き出します。  
MPI環境の終了処理を行います。

```
[aoyama@node01]~/u/aoyama: mpirun -np 3 a.out
MYRANK = 0 IRECV = 1
MYRANK = 1 IRECV = 2
MYRANK = 2 IRECV = 3
```

# 集団通信サブルーチン MPI\_SCATTER

## 機能 (関連する節: 3 - 3 節)

- コミュニケータ(comm)内の、1つの送信元プロセス(root)の送信バッファァー(sendbuf)から、全プロセスの受信バッファァー(recvbuf)にメッセージを送信します。
- MPI\_SCATTERと異なり、送信メッセージの長さ<sup>1</sup>と送信バッファァー(sendbuf)内の位置(変位)を、宛先プロセスごとに変えることができます。

## 使用法

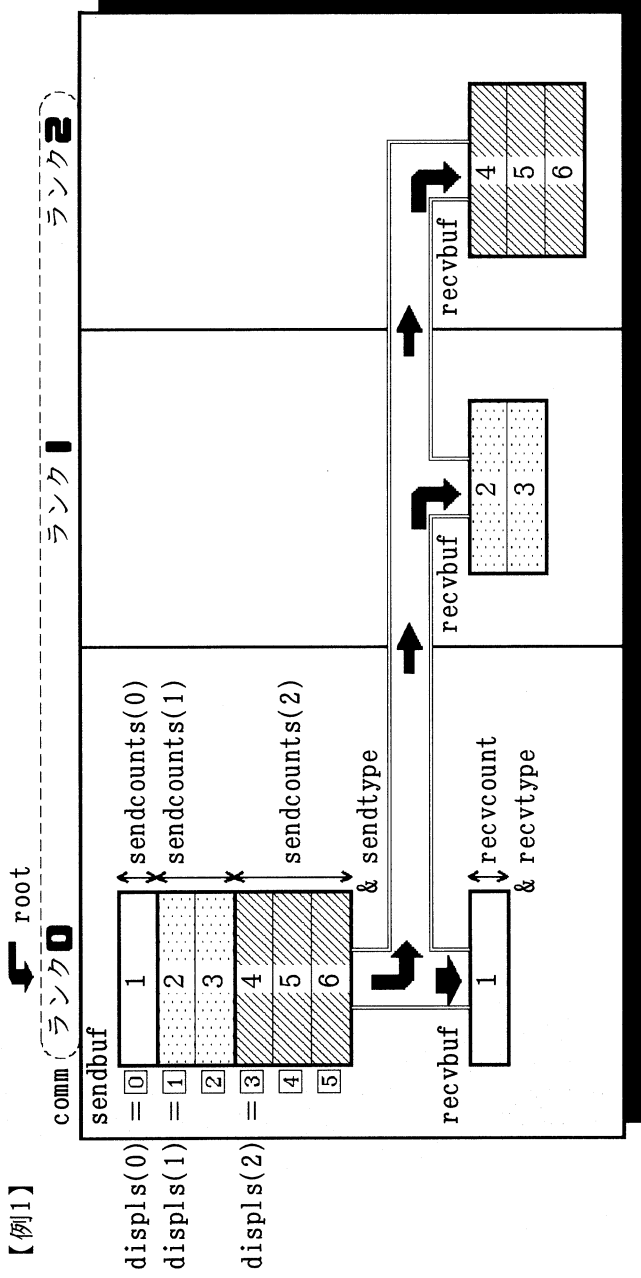
```
CALL MPI_SCATTERV(sendbuf, sendcounts, displs, sendtype,  
recvbuf, recvcount, root, comm, ierror)
```

- sendbuf : 送信バッファァーの先頭アドレスを指定します。rootプロセスのみで意味を持ちます。
- sendcounts : 整数配列。ランクiのプロセスに送信する送信メッセージの要素数を、この配列のi+1番目の要素に指定します(サンプルプログラムに示すように、配列を0から開始した方が分かりやすくなります)。rootプロセスのみで意味を持ちます。
- displs : 整数配列。ランクiのプロセスに送信する送信メッセージが置かれている送信バッファァー sendbuf内の位置(変位; 【例1】参照)を、この配列のi+1番目の要素に指定します(サンプルプログラムに示すように、配列を0から開始した方が分かりやすくなります)。
- sendtype : 整数。送信メッセージのデータを単位として表します。rootプロセスのみで意味を持ちます。
- recvbuf : 受信バッファァーの先頭アドレスを指定します。
- sendbuf : 送信バッファァーと受信バッファァーの実際に使用する部分は、メモリー上で重なってはいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3 - 8 - 1 節を参照)。
- recvcount : 整数。受信メッセージの要素数を指定します。
- recvtype : 整数。受信メッセージのデータを指定します。
- root : 整数。送信元プロセスのcomm内でのランクを指定します。comm内の全プロセスが同じ値を指定する必要があります。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- 各プロセスへの送信メッセージ間でバイト数は異なっても構いません。送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでsendcounts, sendtype, recvcount, recvtypeは異なっても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- 送信元プロセス(root)以外のプロセスでは、sendbuf, sendcounts, displs, sendtypeは無視されますので、適当な変数を指定してもかまいません。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要がありません。

【例1】



■ サンプルプログラム 【例1】 のプログラム例です。

```
PROGRAM MAIN
```

```
INCLUDE 'mpif.h'
```

```
INTEGER ISEND(6), IRECV(3)
```

```
INTEGER ISCNT(0:2), IDISP(0:2)
```

```
DATA ISCNT/1,2,3/ IDISP/0,1,3/
```

```
CALL MPI_INIT(IERR)
```

```
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR) プロセス数NPROCSを取得します。
```

```
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。
```

```
IF (MYRANK == 0) THEN
```

```
DO I = 1, 6
```

```
ISEND(I) = I
```

```
ENDDO
```

```
ENDIF
```

```
CALL MPI_SCATTERV(ISEND,ISCNT ,IDISP,MPI_INTEGER, 通信を行います。
```

```
& IRECV,MYRANK+1 ,MPI_INTEGER,
```

```
& 0,MPI_COMM_WORLD,IERR)
```

```
PRINT *,'MYRANK = ',MYRANK,'IRECV = ',IRECV 通信後の受信バッファの内容を書き出します。
```

```
CALL MPI_FINALIZE(IERR)
```

```
END
```

```
[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
```

```
MYRANK = 0 IRECV = 1 0 0
```

```
MYRANK = 1 IRECV = 2 3 0
```

```
MYRANK = 2 IRECV = 4 5 6
```

MPIが使用するインクルードファイルを指定します。  
送受信バッファを指定します。

送信メッセージの要素数と変位を指定します。

MPI環境の初期化処理を行います。

ランク0のプロセスは送信バッファにデータを  
セットします。

# 集団通信サブルーチン MPI\_GATHER

## 機能 (関連する節: 3-3-5 節)

- コミュニケータ(comm)内の、全プロセスの送信バッファ(sendbuf)から、1つの宛先プロセス(root)の受信バッファ(recvbuf)にメッセージを送信します。
- 各送信元プロセスからの受信メッセージの長さは一定で、受信バッファ(recvbuf)の先頭から送信元プロセスのランクが小さい順に入ります。
- 本サブルーチンはMPI\_SCATTERと逆の通信になります。

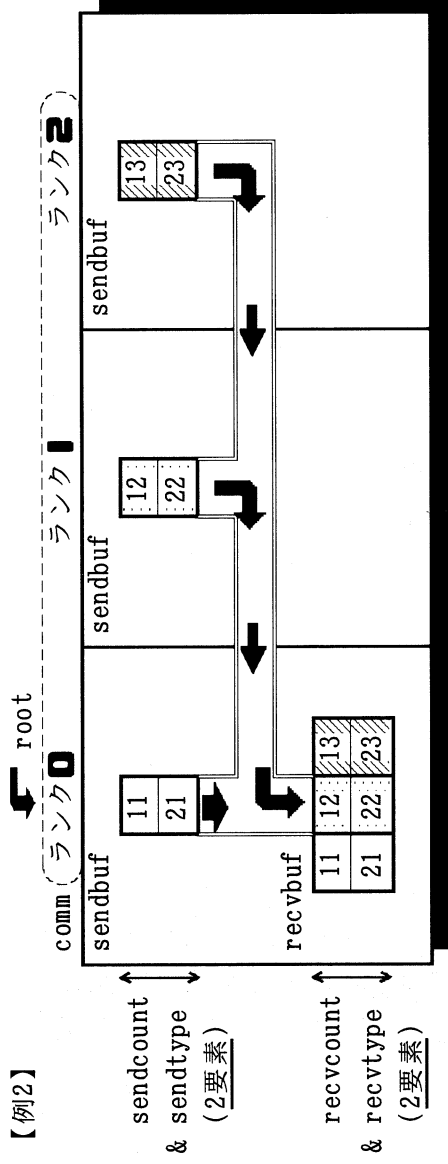
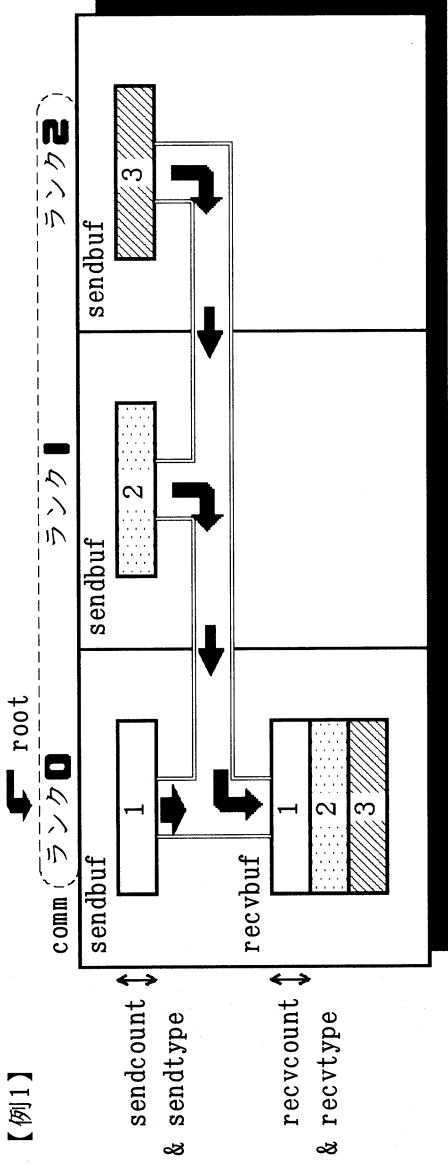
## 使用法

```
CALL MPI_GATHER(sendbuf, sendcount, sendtype,  
recvbuf, recvcount, recvtype, root, comm, ierror)
```

- sendbuf : 送信バッファの先頭アドレスを指定します。
  - sendcount : 整数。送信メッセージの要素数を指定します。
  - sendtype : 整数。送信メッセージのデータ型を指定します。
  - recvbuf : 受信バッファの先頭アドレスを指定します。rootプロセスのみで意味を持ちます。
- 送信バッファと受信バッファの実際使用する部分は、メモリー上で重なってはいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1節を参照)。
- recvcount : 整数。1つのプロセスからの受信メッセージの要素数を指定します。
  - recvtype : 整数。rootプロセスのみで意味を持ちます。
  - root : 整数。宛先プロセスのcomm内でのランクを指定します。
  - comm : 整数。全プロセスが同じ値を指定する必要があります。
  - ierror : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。

## 注意

- 各プロセスからの送信メッセージ間、および、送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでsendcount, sendtype, recvcount, recvtypeは異なっていても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- 宛先プロセス(root)以外のプロセスでは、recvbuf, recvcount, recvtypeは無視されますので、適当な変数を指定してもかまいません。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。



■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER IRECV(3)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)
  ISEND = MYRANK + 1
  CALL MPI_GATHER(ISEND,1,MPI_INTEGER,0,MPI_COMM_WORLD,IERR,
    & IRECV,1,MPI_INTEGER,0,MPI_COMM_WORLD,IERR)
  IF (MYRANK == 0) PRINT *,IRECV = ',IRECV
  CALL MPI_FINALIZE(IERR)
  END
  
```

```

[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
  IRECV = 1 2 3
  
```

# 集団通信サブルーチン MPI\_GATHER

## 機能 (関連する節: 4-6-3-1 節)

- コミュニケータ(comm)内の、全プロセスの送信バッファ(sendbuf)から、1つの宛先プロセス(root)の受信バッファ(recvbuf)にメッセージを送信します。
- MPI\_GATHERと異なり、受信メッセージの長さと受信バッファ(recvbuf)内の位置(変位)を、送信元プロセスごとに変えることができます。

## 使用法

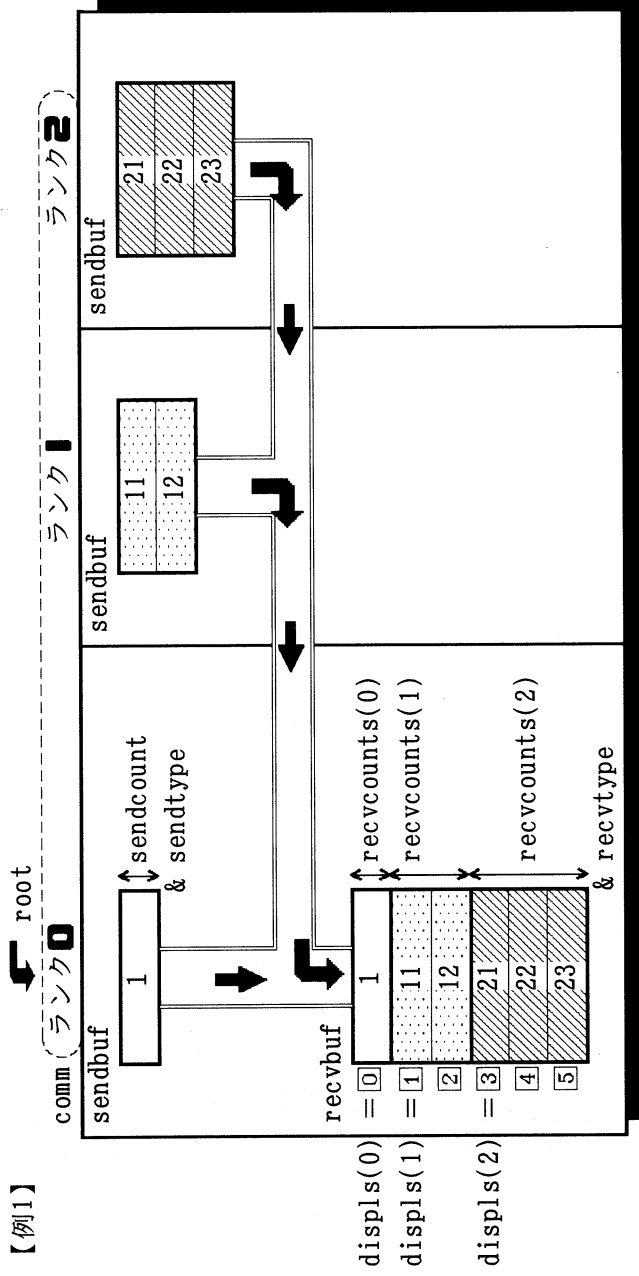
```
CALL MPI_GATHER(sendbuf, sendcount, sendtype,
                recvbuf, recvcnts, displs, root, comm, ierror)
```

- sendbuf : 送信バッファの先頭アドレスを指定します。
- sendcount : 整数。送信メッセージの要素数を指定します。
- sendtype : 整数。送信メッセージのデータ型を指定します。
- recvbuf : 受信バッファの先頭アドレスを指定します。rootプロセスのみで意味を持ちます。送信バッファと受信バッファの実際に使用する部分は、メモリー上で重なってははいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1節を参照)。
- recvcnts : 整数配列。ランクiのプロセスからの受信メッセージの要素数を、この配列のi+1番目の要素に指定します(サブプログラムに示すように、配列を0から開始した方が分かりやすくなります)。rootプロセスのみで意味を持ちます。
- displs : 整数配列。ランクiのプロセスからの受信メッセージが置かれる受信バッファrecvbuf内の位置(変位;【例1】参照)を、この配列のi+1番目の要素に指定します(サブプログラムに示すようにに、配列を0から開始した方が分かりやすくなります)。
- recvtpe : 整数。受信メッセージのデータ型を指定します。rootプロセスのみで意味を持ちます。
- root : 整数。宛先プロセスのrankを指定します。
- comm : 整数。送信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- 各プロセスからの送信メッセージ間でバイト数は異なっても構いません。送信メッセージとそれに対応する受信メッセージ間のバイト数は一致する必要があります。この条件を満足していれば、各プロセスでsendcount, sendtype, recvcnts, recvtpeは異なっていても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- 宛先プロセス(root)以外のプロセスでは、recvbuf, recvcnts, displs, recvtpeは無視されるので、適当な変数を指定してもかまいません。
- 1回のコールで受信バッファの同じ位置に複数の送信元プロセスからメッセージを受信しないで下さい。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。

【例1】



■ サンプルプログラム 【例1】のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER ISEND(3), IRECV(6)
INTEGER IRCNT(0:2), IDISP(0:2)
DATA IRCNT/1,2,3/ IDISP/0,1,3/
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。
DO I = 1, MYRANK+1
  ISEND(I) = I + MYRANK*10
ENDDO
CALL MPI_GATHERV(ISEND,MYRANK+1 ,MPI_INTEGER,
& IRECV,IRCNT ,IDISP,MPI_INTEGER,
& 0,MPI_COMM_WORLD,IERR)
IF (MYRANK == 0) PRINT *, 'IRECV = ', IRECV 通信後の受信バッファの内容を書き出します。
CALL MPI_FINALIZE(IERR)
END
    
```

```

[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
IRECV = 1 11 12 21 22 23
    
```

## 機能 (関連する節: 3-3 節)

- コミュニケータ(comm)内の、全プロセスの送信バッファ(sendbuf)から、全プロセスの受信バッファ(recvbuf)にお互いにメッセージを送信します。
- 各送信元プロセスからの受信メッセージの長さは一定で、受信バッファ(recvbuf)の先頭から送信元プロセスのランクが小さい順に入ります。
- 本サブルーチンは、メッセージをMPI\_GATHERで1つのプロセスに収集した後、それをMPI\_BCASTで全プロセスに送信したのと機能的に等価です。使用方法も『宛先プロセスのランク(root)』の引数がない点を除き、MPI\_GATHERと同じです。

## 使用法

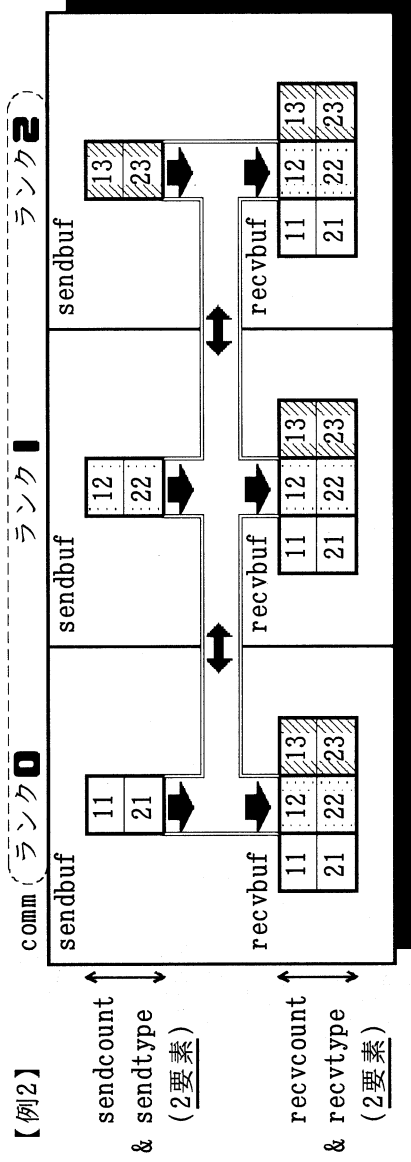
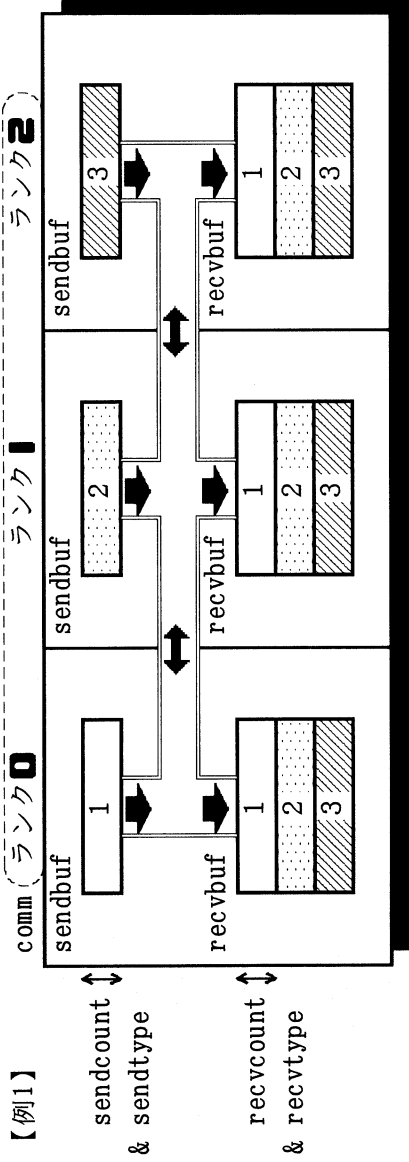
```
CALL MPI_ALLGATHER(sendbuf, sendcount, sendtype,
                  recvbuf, recvcnt, recvtype, comm, ierror)
```

- sendbuf : 送信バッファの先頭アドレスを指定します。
- sendcount : 整数。送信メッセージの要素数を指定します。
- sendtype : 整数。送信メッセージのデータ型を指定します。
- recvbuf : 受信バッファの先頭アドレスを指定します。  
送信バッファと受信バッファの実際使用する部分は、メモリー上で重なってはいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1節を参照)。
- recvcnt : 整数。1つのプロセスからの受信メッセージの要素数を指定します。
- recvtype : 整数。受信メッセージのデータ型を指定します。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケータを指定します。全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- 各プロセスからの送信メッセージ間、および、送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでsendcount, sendtype, recvcnt, recvtypeは異なっていても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要がありません。





■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER IRECV(3)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)  プロセス数NPROCSを取得します。
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)  自分のランクMYRANKを取得します。
  ISEND = MYRANK + 1
  CALL MPI_ALLGATHER(ISEND, 1, MPI_INTEGER,
    &
    IRECV, 1, MPI_INTEGER, MPI_COMM_WORLD, IERR)  通信を行います。
  PRINT *, 'MYRANK = ', MYRANK, 'IRECV = ', IRECV  通信後の受信バッファの内容を書き出します。
  CALL MPI_FINALIZE(IERR)
  END
  
```

```

[aoyama@node01]~/aoyama: mpirun -np 3 a.out
MYRANK = 0 IRECV = 1 2 3
MYRANK = 1 IRECV = 1 2 3
MYRANK = 2 IRECV = 1 2 3
  
```

# 集団通信サブルーチン MPI\_ALLGATHERV

## 機能 (関連する節: 4-6-3-2 節)

- コミュニケータ (comm) 内の、全プロセスの送信バッファ (sendbuf) から、全プロセスの受信バッファ (recvbuf) にお互いにメッセージを送信します。
- MPI\_ALLGATHER と異なり、受信メッセージの長さ (recvbuf) と送信バッファ (sendbuf) 内の位置 (変位) を、送信元プロセスごとに変えることができます。
- 使用法は『宛先プロセスのランク (root)』の引数がない点を除き、MPI\_GATHERV と同じです。

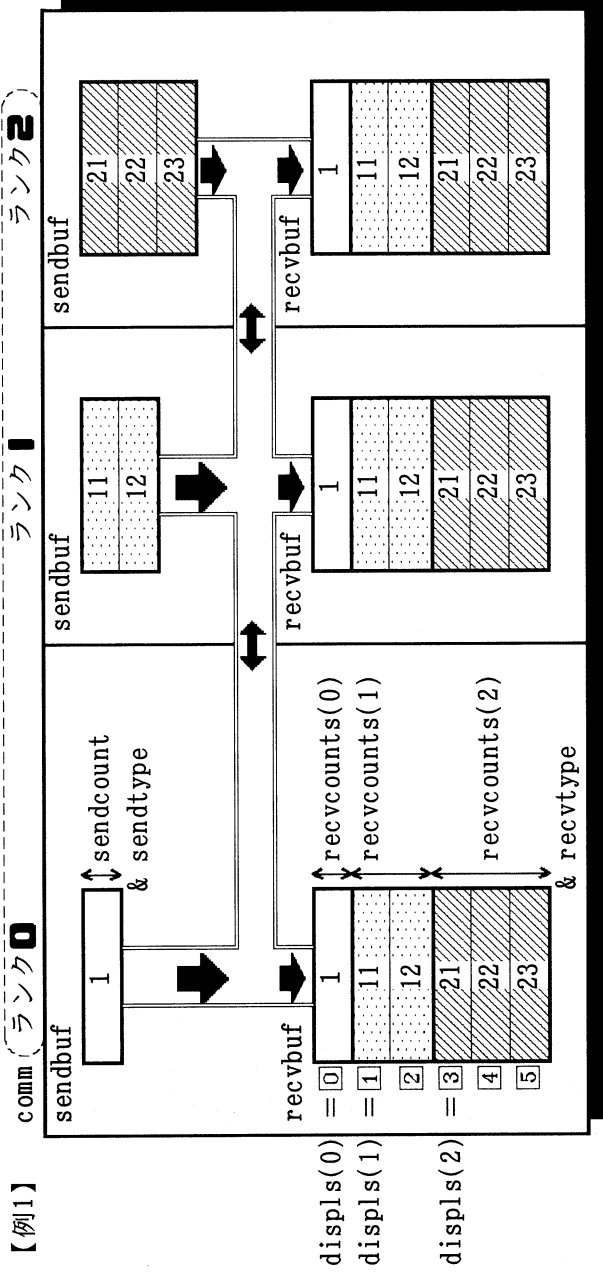
## 使用法

```
CALL MPI_ALLGATHERV(sendbuf, sendcount, sendtype,
                    recvbuf, recvcnts, displs, recvtype, comm, ierror)
```

- sendbuf : 送信バッファの先頭アドレスを指定します。
- sendcount : 整数。送信メッセージの要素数を指定します。
- sendtype : 整数。送信メッセージのデータ型を指定します。
- recvbuf : 受信バッファの先頭アドレスを指定します。
- sendcount, sendtype, recvcount, displ, recvtype, comm, ierror : 送信バッファと受信バッファの実際に使用する部分は、メモリー上で重なってははいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1 節を参照)。
- recvcnts : 整数配列。ランクiのプロセスからの受信メッセージの要素数を、この配列のi+1番目の要素に指定します(サンプルプログラムに示すように、配列を0から開始した方が分かりやすくなります)。
- displs : 整数配列。ランクiのプロセスからの受信メッセージが置かれる受信バッファrecvbuf内の位置(変位; 【例1】参照)を、この配列のi+1番目の要素に指定します(サンプルプログラムに示すように、配列を0から開始した方が分かりやすくなります)。この値はタイプrecvtypeを単位として表します。
- recvtype : 整数。受信メッセージのデータ型を指定します。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- 各プロセスからの送信メッセージ間でバイト数は異なっても構いません。送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでsendcount, sendtype, recvcnts, recvtypeは異なっても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- 1回のコールで受信バッファの同じ位置に複数の送信元プロセスからメッセージを受信しないで下さい。
- コミュニケータ (comm) 内の全プロセスが本サブルーチンをコールする必要がありません。



■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER ISEND(3), IRECV(6)
INTEGER IRCNT(0:2), IDISP(0:2)
DATA IRCNT/1,2,3/ IDISP/0,1,3/
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)
DO I = 1, MYRANK+1
  ISEND(I) = I + MYRANK*10
ENDDO
CALL MPI_ALLGATHERV(ISEND,MYRANK+1, MPI_INTEGER,
& IRECV,IRCNT, IDISP,MPI_INTEGER,
& MPI_COMM_WORLD,IERR)
PRINT *, 'MYRANK = ',MYRANK, 'IRECV = ',IRECV
CALL MPI_FINALIZE(IERR)
END

```

MPIが使用するインクルードファイルを指定します。  
送受信バッファを指定します。

受信メッセージの要素数と変位を指定します。

MPI環境の初期化処理を行います。

CALL MPI\_COMM\_SIZE(MPI\_COMM\_WORLD,NPROCS,IERR) プロセス数NPROCSを取得します。

CALL MPI\_COMM\_RANK(MPI\_COMM\_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。

各プロセスは送信バッファにデータをセットします。

CALL MPI\_ALLGATHERV(ISEND,MYRANK+1, MPI\_INTEGER, 通信を行います。  
& IRECV,IRCNT, IDISP,MPI\_INTEGER,  
& MPI\_COMM\_WORLD,IERR)

PRINT \*, 'MYRANK = ',MYRANK, 'IRECV = ',IRECV 通信後の受信バッファの内容を書き出します。  
CALL MPI\_FINALIZE(IERR) MPI環境の終了処理を行います。

```

[aoyama@node01]~/aoyama: mpirun -np 3 a.out
MYRANK = 0 IRECV = 1 11 12 21 22 23
MYRANK = 1 IRECV = 1 11 12 21 22 23
MYRANK = 2 IRECV = 1 11 12 21 22 23

```

## 機能 (関連する節: 4-6-10 節)

- コミュニケータ (comm) 内の、全プロセスの送信バッファ (sendbuf) から、全プロセスの受信バッファ (recvbuf) にお互いにメッセージを送信します。
- 各宛先プロセスへの送信メッセージの長さは一定で、送信バッファ (sendbuf) の先頭から宛先プロセスのランクが小さい順に送信されます。
- 各送信元プロセスからの受信メッセージの長さは一定で、受信バッファ (recvbuf) の先頭から送信元プロセスのランクが小さい順に入ります。
- 送信側に着目すると MPI\_SCATTER を行い (例えば【例1】で、ランク 0 の送信バッファに入っている 1, 2, 3 が各プロセスに送られます)、受信側に着目すると MPI\_GATHER を行っています (例えば【例1】で、各プロセスから送られた 1, 11, 21 がランク 0 の受信バッファに入ります)。

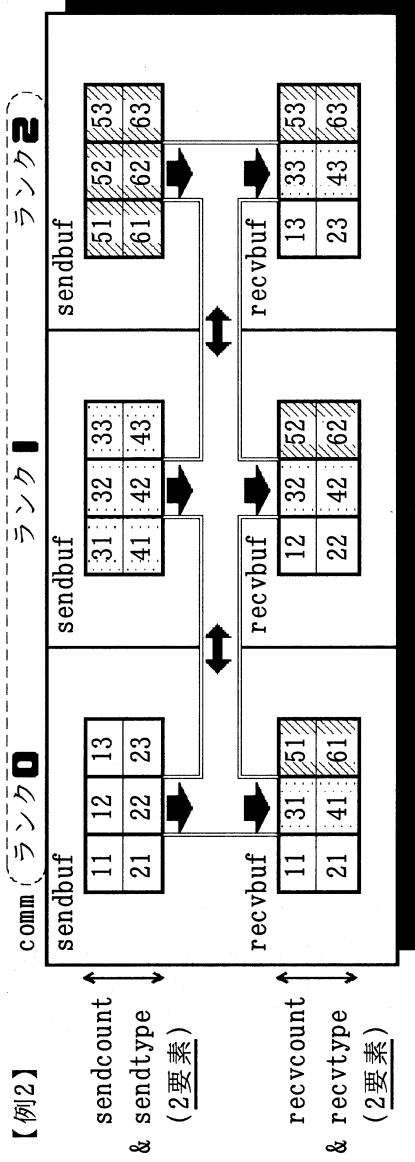
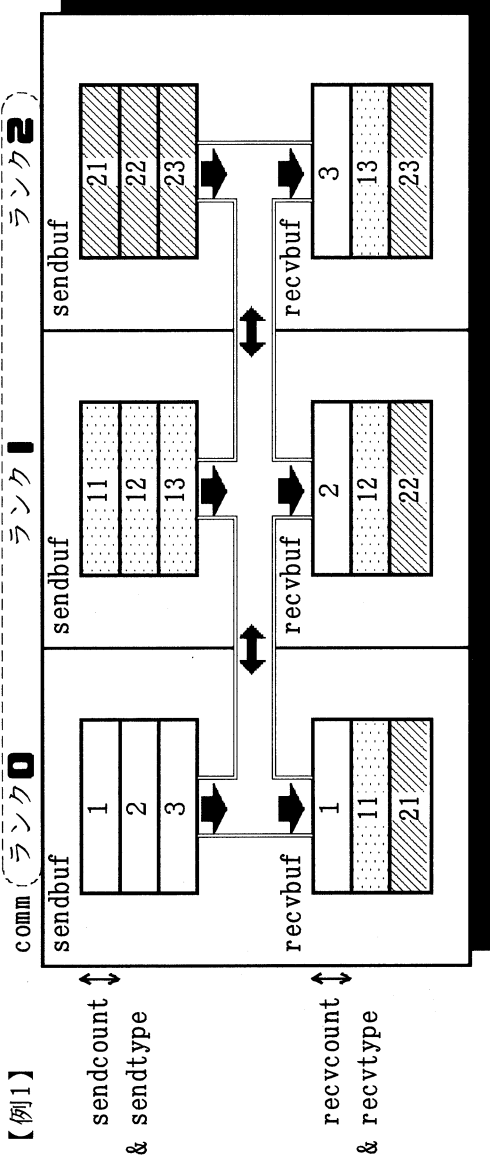
## 使用法

```
CALL MPI_ALLTOALL (sendbuf, sendcount, sendtype,
                  recvbuf, recvcount, recvtype, comm, ierror)
```

- **sendbuf** : 送信バッファの先頭アドレスを指定します。
- **sendcount** : 整数。1つのプロセスに送信する送信メッセージの要素数を指定します。
- **sendtype** : 整数。送信メッセージのデータ型を指定します。
- **recvbuf** : 受信バッファの先頭アドレスを指定します。  
送信バッファと受信バッファの実際使用部分は、メモリー上で重なってはいいけません (MPI-2 の MPI\_IN\_PLACE を使用することはできません)。
- **recvcount** : 整数。1つのプロセスからの受信メッセージの要素数を指定します。
- **recvtype** : 整数。受信メッセージのデータ型を指定します。
- **comm** : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。全プロセスが同じ値を指定する必要があります。
- **ierror** : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- 各プロセスへの送信メッセージ間、および、送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスで sendcount, sendtype, recvcount, recvtype は異なっても構いません (ただしデータ型 (INTEGER, REAL など) は一致している必要があります)。
- コミュニケータ (comm) 内の全プロセスが本サブルーチンをコールする必要がありません。



■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISEND(3), IRECV(3)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  DO I = 1, NPROCS
    ISEND(I) = I + MYRANK*10
  ENDDO
  CALL MPI_ALLTOALL(ISEND, 1, MPI_INTEGER,
    & IRECV, 1, MPI_INTEGER, MPI_COMM_WORLD, IERR)
  PRINT *, 'MYRANK = ', MYRANK, 'IRECV = ', IRECV
  CALL MPI_FINALIZE(IERR)
END

```

```

[aoyama@node01]~/aoyama: mpirun -np 3 a.out
MYRANK = 0 IRECV = 1 11 21
MYRANK = 1 IRECV = 2 12 22
MYRANK = 2 IRECV = 3 13 23

```

MPIが使用するインクルードファイルを指定します。  
送受信バッファを指定します。

MPI環境の初期化処理を行います。  
プロセス数NPROCSを取得します。  
自分のランクMYRANKを取得します。

各プロセスは送信バッファにデータをセットします。

通信を行います。

通信後の受信バッファの内容を書き出します。  
MPI環境の終了処理を行います。

**機能** (関連する節: 4-6-3-3節, 4-6-1-0節, 5-2-3節)

- コミュニケータ(comm)内の、全プロセスの送信バッファアー(sendbuf)から、全プロセスの受信バッファアー(recvbuf)にお互いにメッセージを送信します。
- MPI\_ALLTOALLとは異なり、送信メッセージの長さ(=sendbuf)と送信バッファアー(sendbuf)内の位置(変位)を、宛先プロセスごとに変えることができます。
- MPI\_ALLTOALLとは異なり、受信メッセージの長さ(=recvbuf)と受信バッファアー(recvbuf)内の位置(変位)を、送信元プロセスごとに変えることができます。
- 【例1】では、全プロセスの送信バッファアーが同じ構造になっていますが、異なっても構いません。

**使用法**

```
CALL MPI_ALLTOALLV(sendbuf, sendcounts, sdispls, sendtype,  
recvbuf, recvcounts, rdispls, recvtype, comm, ierror)
```

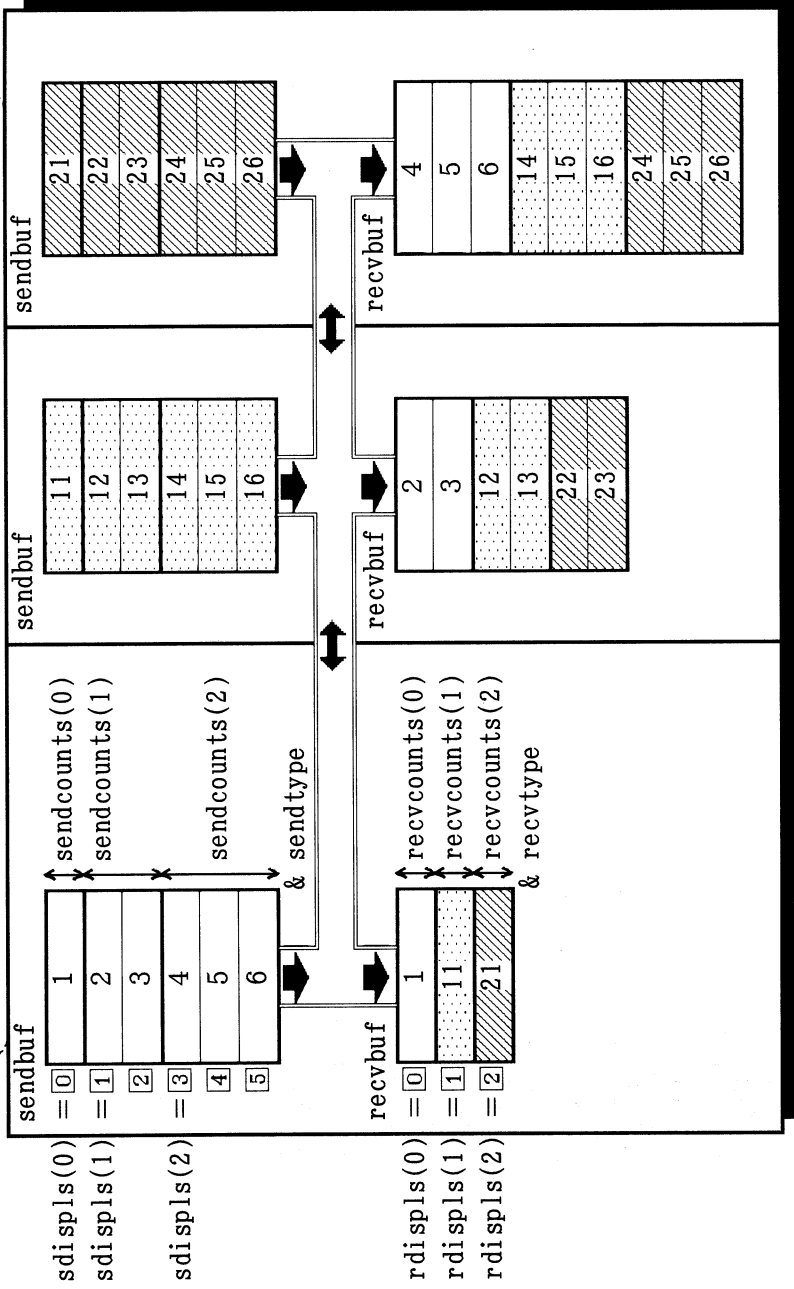
- sendbuf : 送信バッファアーの先頭アドレスを指定します。
- sendcounts : 整数配列。ランクiのプロセスに送信する送信メッセージの要素数を、この配列のi+1番目の要素に指定します(サンプリングプログラムに示すように、配列を0から開始した方が分かりやすくなります)。
- sdispls : 整数配列。ランクiのプロセスに送信する送信メッセージが置かれている送信バッファアーrecvbuf内の位置(変位;【例1】参照)を、この配列のi+1番目の要素に指定します(サンプリングプログラムに示すように、配列を0から開始した方が分かりやすくなります)。
- sendtype : この値はタイプsendtypeを単位として表します。
- recvbuf : 整数。送信メッセージのデータ型を指定します。
- recvbuf : 受信バッファアーの先頭アドレスを指定します。
- sendbuf : 送信バッファアーと受信バッファアーの実際に使用する部分は、メモリー上で重なってはいけません(MPI-2のMPI\_IN\_PLACEを使用することはできません)。
- recvcounts : 整数配列。ランクiのプロセスからの受信メッセージの要素数を、この配列のi+1番目の要素に指定します(サンプリングプログラムに示すように、配列を0から開始した方が分かりやすくなります)。
- rdispls : 整数配列。ランクiのプロセスからの受信メッセージが置かれる受信バッファアーrecvbuf内の位置(変位;【例1】参照)を、この配列のi+1番目の要素に指定します(サンプリングプログラムに示すように、配列を0から開始した方が分かりやすくなります)。
- recvtype : この値はタイプrecvtypeを単位として表します。
- comm : 整数。受信メッセージのデータ型を指定します。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

**注意**

- 各プロセスへの送信メッセージ間でバイト数は異なっても構いません。送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでsendcounts, sendtype, recvcounts, recvtypeは異なっても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- 1回のコールで受信バッファアーの同じ位置に複数の送信元プロセスからメッセージを受信しないで下さい。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。

【例1】

ランク 0      ランク 1      ランク 2



■ サンプルプログラム 【例1】 のプログラム例です。

PROGRAM MAIN

INCLUDE 'mpif.h'

INTEGER ISEND(6), IRECV(9)

INTEGER ISCNT(0:2), ISDSP(0:2), IRCNT(0:2), IRDSP(0:2)

DATA ISCNT/1,2,3/ ISDSP/0,1,3/

CALL MPI\_INIT(IERR)

CALL MPI\_COMM\_SIZE(MPI\_COMM\_WORLD, NPROCS, IERR)      プロセス数NPROCSを取得します。  
 CALL MPI\_COMM\_RANK(MPI\_COMM\_WORLD, MYRANK, IERR)      自分のランクMYRANKを取得します。

DO I = 0, NPROCS-1

    IRCNT(I) = 1 + MYRANK

    IRDSP(I) = I\*(MYRANK+1)

ENDDO

DO I = 1, 6

    ISEND(I) = I + MYRANK\*10

ENDDO

CALL MPI\_ALLTOALLV(ISEND, ISCNT, ISDSP, MPI\_INTEGER, IRECV, IRCNT, IRDSP, MPI\_INTEGER,

&      MPI\_COMM\_WORLD, IERR)

&      MPI\_COMM\_WORLD, IERR)

PRINT \*, 'MYRANK = ', MYRANK, 'IRECV = ', IRECV      通信後の受信バッファの内容を書き出します。  
 CALL MPI\_FINALIZE(IERR)      MPI環境の終了処理を行います。  
 END

[aoyama@node01]/u/aoyama: mpirun -np 3 a.out

MYRANK = 0 IRECV = 1 11 21 0 0 0 0

MYRANK = 1 IRECV = 2 3 12 13 22 23 0 0

MYRANK = 2 IRECV = 4 5 6 14 15 16 24 25 26

MPIが使用するインクルードファイルを指定します。  
 送受信バッファを指定します。

送信メッセージの要素数と変位を指定します。

MPI環境の初期化処理を行います。

受信メッセージの要素数を指定します。

受信メッセージの変位を指定します。

各プロセスは送信バッファにデータをセットします。

### 機能 (関連する節: 4-6-3-3 節)

- サブルーチンはMPI\_ALLTOALLVとほぼ同じで、「使用法」の波線の部分だけが異なります。
- コミュニケータ(comm)内の、全プロセスの送信バッファァー(sendbuf)から、全プロセスの受信バッファァー(recvbuf)にお互いにメッセージを送信します。
- MPI\_ALLTOALLとは異なり、送信メッセージの長さ、送信バッファァー(sendbuf)内の位置(変位)、およびデータ型を、宛先プロセスごとに変えることができます。
- MPI\_ALLTOALLとは異なり、受信メッセージの長さ、受信バッファァー(recvbuf)内の位置(変位)、およびデータ型を、送信元プロセスごとに変えることができます。
- 【例1】では、全プロセスの送信バッファァーが同じ構造になっていますが、異なっても構いません。

### 使用法

```
CALL MPI_ALLTOALLW(sendbuf, sendcounts, sdispls, sendtypes,
recvbuf, recvcounts, rdispls, recvtypes, comm, ierror)
```

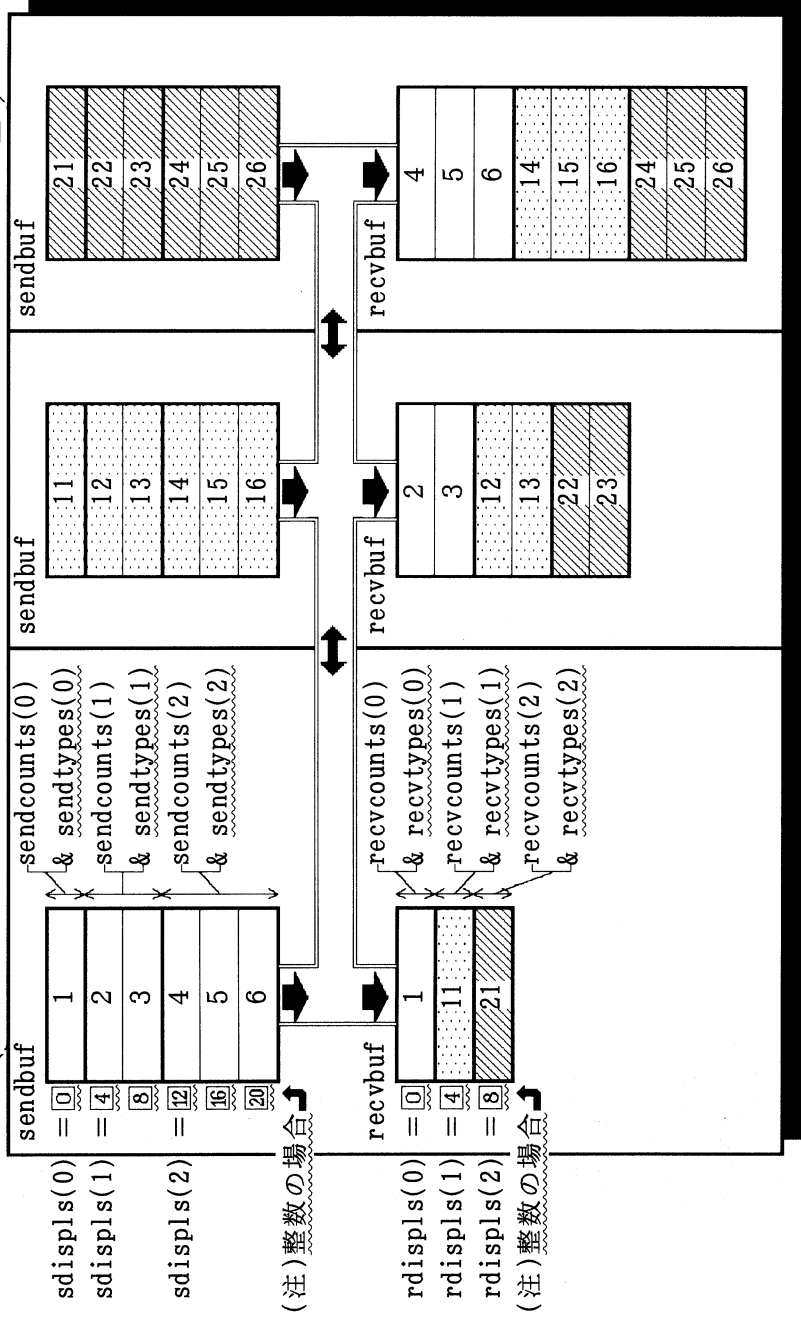
- **sendbuf** : 送信バッファァーの先頭アドレスを指定します。
- **sendcounts** : 整数配列。ランクiのプロセスに送信する送信メッセージの要素数を、この配列のi+1番目の要素に指定します(注1を参照して下さい)。
- **sdispls** : 整数配列。ランクiのプロセスに送信する送信メッセージが置かれている送信バッファァーsendbuf内の位置(変位;【例1】参照)を、この配列のi+1番目の要素に指定します(注1を参照して下さい)。
- **sendtypes** : 整数配列。ランクiのプロセスに送信する送信メッセージのデータ型を、この配列のi+1番目の要素に指定します(注1を参照して下さい)。
- **recvbuf** : 受信バッファァーの先頭アドレスを指定します。  
送信バッファァーと受信バッファァーの実際に使用する部分は、メモリー上で重なってはいけません(MPI-2のMPI\_IN\_PLACEを使用することはできません)。
- **recvcounts** : 整数配列。ランクiのプロセスからの受信メッセージの要素数を、この配列のi+1番目の要素に指定します(注1を参照して下さい)。
- **rdispls** : 整数配列。ランクiのプロセスからの受信メッセージが置かれる受信バッファァーrecvbuf内の位置(変位;【例1】参照)を、この配列のi+1番目の要素に指定します(注1を参照して下さい)。
- **recvtypes** : 整数配列。ランクiのプロセスからの受信メッセージのデータ型を、この配列のi+1番目の要素に指定します(注1を参照して下さい)。
- **comm** : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。全プロセスが同じ値を指定する必要があります(注1を参照して下さい)。
- **ierror** : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

### 注意

- (注1) サンプルプログラムに示すように、配列を0から開始した方が分かりやすくなります。
- 各プロセスへの送信メッセージ間でバイト数は異なっても構いません。送信メッセージとそれに対応する受信メッセージ間のバイト数は一致している必要があります。この条件を満足していれば、各プロセスでsendcounts, sendtypes, recvcounts, recvtypesは異なっても構いません(ただしデータ型(INTEGER, REALなど)は一致している必要があります)。
- 1回のコールで受信バッファァーの同じ位置に複数の送信元プロセスからメッセージを受信しないで下さい。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。



【例1】 comm (ランク0) ランク1 ランク2



■ サンプルプログラム 【例1】 のプログラム例です。MPI\_ALLTOALLVの例と異なる部分を波線で示します。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISEND(6), IRECV(9)
  INTEGER ISCNT(0:2), ISDSP(0:2), IRCNT(0:2), IRDSP(0:2)
  INTEGER ISTYPE(0:2), IRTYPE(0:2)
  DATA ISCNT/1,2,3/ ISDSP/0,4,12/
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  DO I = 0, NPROCS-1
    IRCNT(I) = 1 + MYRANK
    IRDSP(I) = I*(MYRANK+1)*4
    ISTYPE(I) = MPI_INTEGER
    IRTYPE(I) = MPI_INTEGER
  ENDDO
  DO I = 1, 6
    ISEND(I) = I + MYRANK*10
  ENDDO
  CALL MPI_ALLTOALLV(ISEND, ISCNT, ISDSP, ISTYPE, IRECV, IRCNT, IRDSP, IRTYPE, MPI_COMM_WORLD, IERR)
  PRINT *, 'MYRANK = ', MYRANK, 'IRECV = ', IRECV
  CALL MPI_FINALIZE(IERR)
  END
  
```

MPIが使用するインクルードファイルを指定します。  
送受信バッファを指定します。

送信メッセージの要素数と変位を指定します。

MPI環境の初期化処理を行います。

CALL MPI\_COMM\_SIZE(MPI\_COMM\_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。

CALL MPI\_COMM\_RANK(MPI\_COMM\_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。

受信メッセージの要素数を指定します。

受信メッセージの変位を指定します。

送信メッセージのデータ型を指定します。

受信メッセージのデータ型を指定します。

各プロセスは送信バッファにデータをセットします。

CALL MPI\_ALLTOALLV(ISEND, ISCNT, ISDSP, ISTYPE, IRECV, IRCNT, IRDSP, IRTYPE, MPI\_COMM\_WORLD, IERR)

PRINT \*, 'MYRANK = ', MYRANK, 'IRECV = ', IRECV 通信後の受信バッファの内容を書き出します。

CALL MPI\_FINALIZE(IERR) MPI環境の終了処理を行います。

END

```

[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
MYRANK = 0 IRECV = 1 11 21 0 0 0 0
MYRANK = 1 IRECV = 2 3 12 13 22 23 0 0
MYRANK = 2 IRECV = 4 5 6 14 15 16 24 25 26
  
```

## 機能 (関連する節: 3-3-6節, 4-6-4節, 4-6-6-1節)

- コミュニケータ(comm)内の、全プロセスの送信バッファァー(sendbuf)のメッセージが、通信しながら演算(op)され、結果が1つの宛先プロセス(root)の受信バッファァー(recvbuf)に入ります。
- 送信バッファァーが配列の場合は、配列の対応する要素ごとに演算が行われます。
- MPIの基本データ型(datatype)とMPIの定義済み演算(op)の組み合わせを以下に示します。演算(op)をユーザーが自分で定義することもできます。

MPIの基本データ型	MPIの定義済み演算
MPI_INTEGER (注1)	MPI_SUM (合計)
MPI_REAL (注2)	MPI_PROD (積)
MPI_DOUBLE_PRECISION (注3)	
MPI_COMPLEX (注4)	
MPI_INTEGER (注1)	MPI_MAX (最大)
MPI_REAL (注2)	MPI_MIN (最小)
MPI_DOUBLE_PRECISION (注3)	
MPI_2INTEGER	MPI_MAXLOC (最大と位置)
MPI_2REAL	MPI_MINLOC (最小と位置)
MPI_2DOUBLE_PRECISION (注5)	

MPIの基本データ型	MPIの定義済み演算
MPI_LOGICAL	MPI_LAND (論理AND)
	MPI_LOR (論理OR)
	MPI_LXOR (論理XOR)
MPI_INTEGER(注1)	MPI_BAND (ビットAND)
MPI_BYTE	MPI_BOR (ビットOR)
	MPI_BXOR (ビットXOR)

下記が使用できるマシンの環境もあります。

- (注1) MPI\_INTEGER4
- (注2) MPI\_REAL4
- (注3) MPI\_REAL8
- (注4) MPI\_COMPLEX16, MPI\_DOUBLE\_COMPLEX
- (注5) MPI\_2REAL8

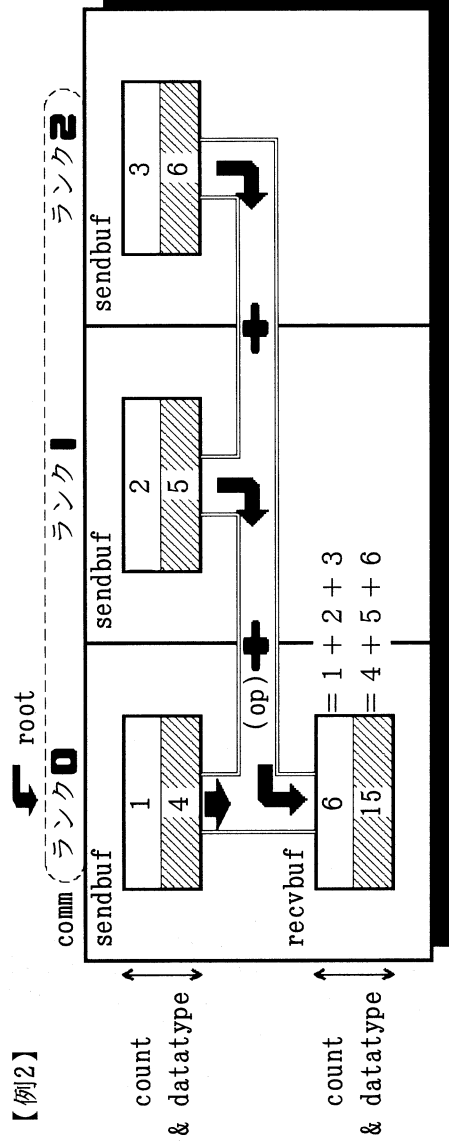
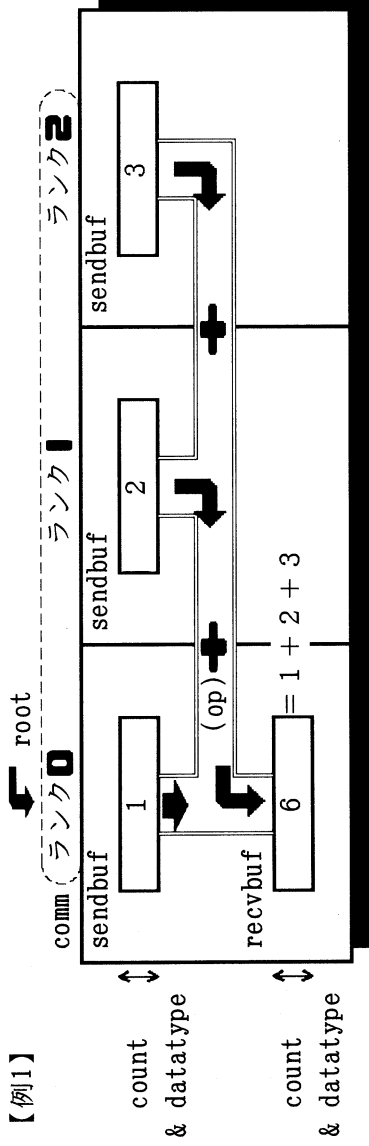
## 使用法

```
CALL MPI_REDUCE(sendbuf, recvbuf, count, datatype, op, root, comm, ierror)
```

- sendbuf : 送信バッファァーの先頭アドレスを指定します。
- recvbuf : 受信バッファァーの先頭アドレスを指定します。rootプロセスのみで意味を持ちます。  
送信バッファァーと受信バッファァーの実際に使用する部分は、メモリー上で重なってははいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1節を参照)。
- count : 整数。送(受)メッセージの要素数を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- datatype : 整数。送(受)メッセージのデータ型を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- op : 整数。演算の種類を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- root : 整数。宛先プロセスのcomm内でのランクを指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケータを指定します。  
全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- 宛先プロセス(root)以外のプロセスでは、recvbufは無視されますので、適当な変数を指定してもかまいません。
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。



■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
ISEND = MYRANK + 1
CALL MPI_REDUCE(ISEND, IRECV, 1, MPI_INTEGER,
&
MPI_SUM, 0, MPI_COMM_WORLD, IERR)
IF (MYRANK == 0) PRINT *, 'IRECV = ', IRECV
CALL MPI_FINALIZE(IERR)
END

```

```

[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
IRECV = 6

```

■ 機能 (関連する節: 4-6-4節, 4-6-6-1節, 5-8節)

- コミュニケータ(comm)内の、全プロセスの送信バッファァー(sendbuf)のメッセージが、通信しながら演算(op)され、結果が全プロセスの受信バッファァー(recvbuf)に入ります。
- 送信バッファァーが配列の場合は、配列の対応する要素ごとに演算が行われます。
- 本サブルーチンは、メッセージをMPI\_REDUCEで1つのプロセスに収集した後、それをMPI\_BCASTで全プロセスに送信したのと機能的に等価です。使用方法も『宛先プロセスのランク(root)』の引数がない点を除き、MPI\_REDUCEと同じです。
- MPIの基本データ型(datatype)とMPIの定義済み演算(op)の組み合わせを以下に示します。演算(op)をユーザが自分で定義することもできます。

MPIの基本データ型	MPIの定義済み演算
MPI_INTEGER	(注1) MPI_SUM (合計)
MPI_REAL	(注2) MPI_PROD (積)
MPI_DOUBLE_PRECISION	(注3)
MPI_COMPLEX	(注4)
MPI_INTEGER	(注1) MPI_MAX (最大)
MPI_REAL	(注2) MPI_MIN (最小)
MPI_DOUBLE_PRECISION	(注3)
MPI_2INTEGER	MPI_MAXLOC(最大と位置)
MPI_2REAL	MPI_MINLOC(最小と位置)
MPI_2DOUBLE_PRECISION	(注5)

MPIの基本データ型	MPIの定義済み演算
MPI_LOGICAL	MPI_LAND(論理AND)
	MPI_LOR (論理OR)
	MPI_LXOR(論理XOR)
MPI_INTEGER(注1)	MPI_BAND(ビットAND)
MPI_BYTE	MPI_BOR (ビットOR)
	MPI_BXOR(ビットXOR)

下記が使用できるマシンの環境もあります。

- (注1) MPI\_INTEGER4
- (注2) MPI\_REAL4
- (注3) MPI\_REAL8
- (注4) MPI\_COMPLEX16, MPI\_DOUBLE\_COMPLEX
- (注5) MPI\_2REAL8

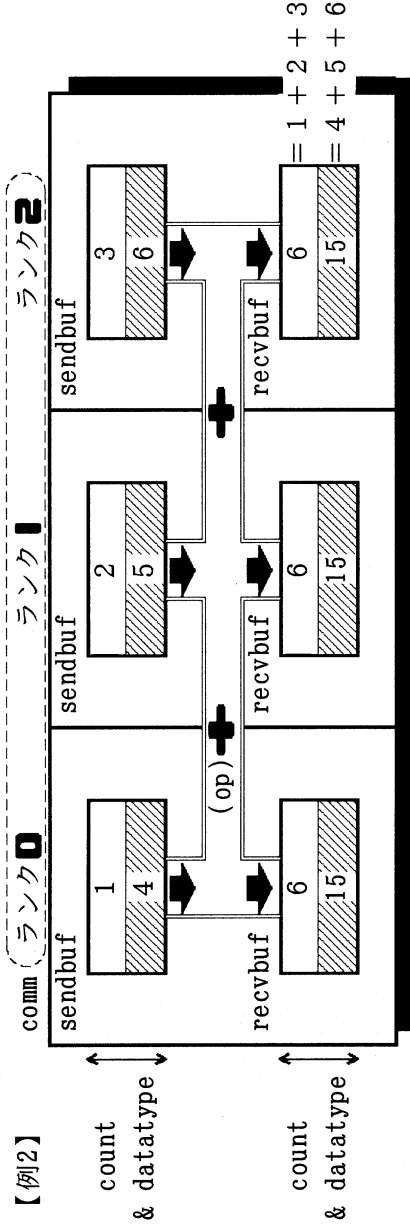
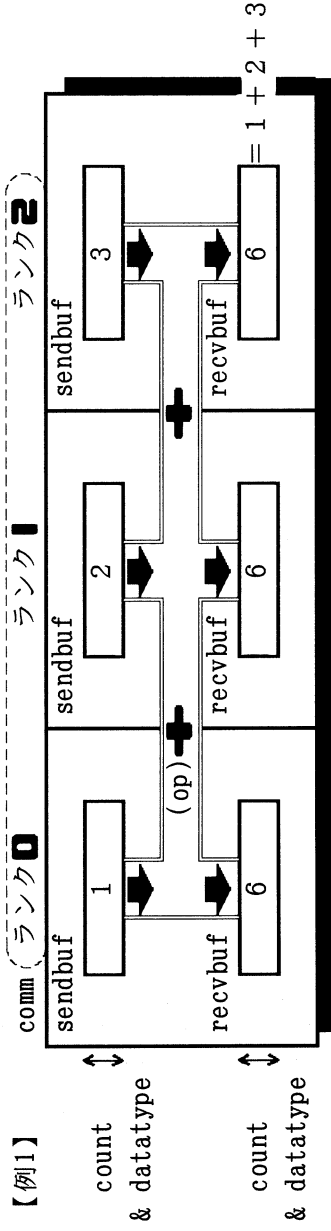
■ 使用法

CALL MPI\_ALLREDUCE(sendbuf, recvbuf, count, datatype, op, comm, ierror)

- sendbuf : 送信バッファァーの先頭アドレスを指定します。
- recvbuf : 受信バッファァーの先頭アドレスを指定します。  
送信バッファァーと受信バッファァーの実際に使用する部分は、メモリー上で重なってはいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1節を参照)。
- count : 整数。送(受)メッセージの要素数を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- datatype : 整数。送(受)メッセージのデータ型を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- op : 整数。演算の種類を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

■ 注意

- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。



■ サンプルプログラム 【例1】 のプログラム例です。

```
PROGRAM MAIN
```

```
INCLUDE 'mpif.h'
```

```
CALL MPI_INIT(IERR)
```

```
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR) プロセス数NPROCSを取得します。
```

```
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。
```

```
ISEND = MYRANK + 1      各プロセスは送信バッファにデータを設定します。
```

```
CALL MPI_ALLREDUCE(ISEND,IRECV,1,MPI_INTEGER,
```

```
& MPI_SUM,MPI_COMM_WORLD,IERR) 通信を行います。
```

```
PRINT *, 'MYRANK = ',MYRANK, 'IRECV = ',IRECV 通信後の受信バッファの内容を書き出します。
```

```
CALL MPI_FINALIZE(IERR)      MPI環境の終了処理を行います。
```

```
END
```

```
[aoyama@node01]~/aoyama: mpirun -np 3 a.out
```

```
MYRANK = 0 IRECV = 6
```

```
MYRANK = 1 IRECV = 6
```

```
MYRANK = 2 IRECV = 6
```

# 集団通信サブルーチン MPI\_SCAN

## 機能 (関連する節: 4-6-9 節)

- コミュニケータ(comm)内の、ランク0~ランクiまでの各プロセスの送信バッファァー(sendbuf)のメッセージが、通信しながら演算(op)され、結果がランクiのプロセスの受信バッファァー(recvbuf)に入ります。これがコミュニティ(comm)内の全プロセスに対して行われます。
- 送信バッファァーが配列の場合は、配列の対応する要素ごとに演算が行われます。
- MPIの基本データ型(datatype)とMPIの定義済み演算(op)の組み合わせを以下に示します。演算(op)をユーザーが自分で定義することもできます。

MPIの基本データ型	MPIの定義済み演算
MPI_INTEGER (注1)	MPI_SUM (合計)
MPI_REAL (注2)	MPI_PROD (積)
MPI_DOUBLE_PRECISION (注3)	
MPI_COMPLEX (注4)	
MPI_INTEGER (注1)	MPI_MAX (最大)
MPI_REAL (注2)	MPI_MIN (最小)
MPI_DOUBLE_PRECISION (注3)	
MPI_2INTEGER	MPI_MAXLOC (最大と位置)
MPI_2REAL	MPI_MINLOC (最小と位置)
MPI_2DOUBLE_PRECISION(注5)	

MPIの基本データ型	MPIの定義済み演算
MPI_LOGICAL	MPI_LAND (論理AND) MPI_LOR (論理OR) MPI_LXOR (論理XOR)
MPI_INTEGER(注1) MPI_BYTE	MPI_BAND (ビットAND) MPI_BOR (ビットOR) MPI_BXOR (ビットXOR)

下記が使用できるマシン環境もあります。

- (注1) MPI\_INTEGER4
- (注2) MPI\_REAL4
- (注3) MPI\_REAL8
- (注4) MPI\_COMPLEX16, MPI\_DOUBLE\_COMPLEX
- (注5) MPI\_2REAL8

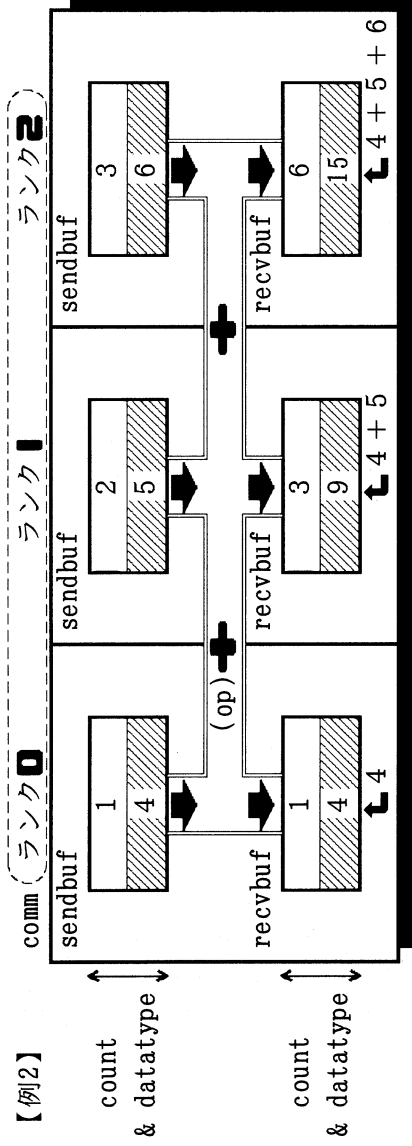
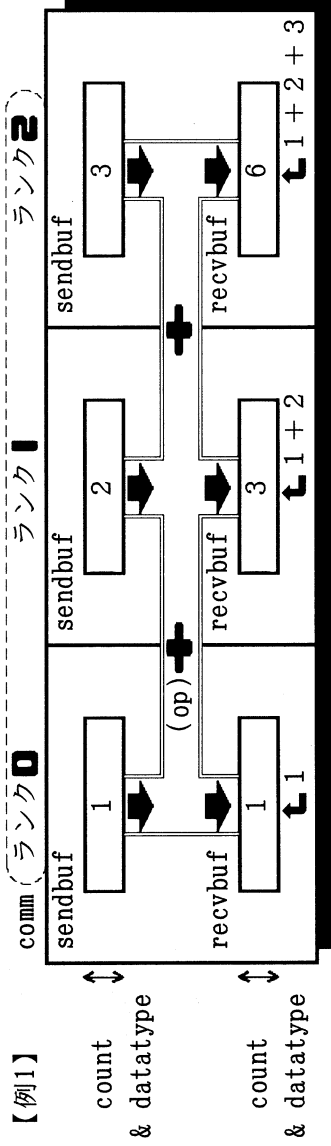
## 使用法

```
CALL MPI_SCAN(sendbuf, recvbuf, count, datatype, op, comm, ierror)
```

- sendbuf : 送信バッファァーの先頭アドレスを指定します。
- recvbuf : 受信バッファァーの先頭アドレスを指定します。  
送信バッファァーと受信バッファァーの実際に使用する部分は、メモリー上で重なってはいけません(ただしMPI-2の新機能MPI\_IN\_PLACEで回避可能です。詳細は3-8-1節を参照)。
- count : 整数。送(受)メッセージの要素数を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- datatype : 整数。送(受)メッセージのデータ型を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。  
整数。演算の種類を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケータを指定します。  
全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要がある場合があります。



■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)
  ISEND = MYRANK + 1
  CALL MPI_SCAN(ISEND,Irecv,1,MPI_INTEGER,
    &
    MPI_SUM,MPI_COMM_WORLD,IERR)
  PRINT *, 'MYRANK = ',MYRANK, 'IRECV = ',IRECV
  CALL MPI_FINALIZE(IERR)
  END
  
```

MPIが使用するインクルードファイルを指定します。  
 MPI環境の初期化処理を行います。  
 プロセス数NPROCSを取得します。  
 自分のランクMYRANKを取得します。  
 各プロセスは送信バッファにデータをセットします。

```

[aoiyama@node01]~/aoiyama: mpirun -np 3 a.out
MYRANK = 0 IRECV = 1
MYRANK = 1 IRECV = 3
MYRANK = 2 IRECV = 6
  
```

## 【MPI-2】集団通信サブルーチン MPI\_EXSCAN

### 機能 (関連する節: 4-6-9 節)

- 通信に参加するプロセスのランクが0~kとします。コミュニケータ(comm)内の、ランク0~ランク*i-1*までの各プロセスの送信バッファァー(sendbuf)のメッセージが、通信しながら演算(op)され、結果がランク*i*のプロセスの受信バッファァー(recvbuf)に入ります。これがコミュニケータ(comm)内のランク1~ランクkのプロセスに対して行われます。ランクkのプロセスの送信バッファァーの値は使用されません。またランク0のプロセスの受信バッファァーに値は設定されません。
- 送信バッファァーが配列の場合は、配列の対応する要素ごとに演算が行われます。
- MPIの基本データ型(datatype)とMPIの定義済み演算(op)の組み合わせを以下に示します。演算(op)をユーザーが自分で定義することもできます。

MPIの基本データ型	MPIの定義済み演算
MPI_INTEGER (注1)	MPI_SUM (合計)
MPI_REAL (注2)	MPI_PROD (積)
MPI_DOUBLE_PRECISION (注3)	
MPI_COMPLEX (注4)	
MPI_INTEGER (注1)	MPI_MAX (最大)
MPI_REAL (注2)	MPI_MIN (最小)
MPI_DOUBLE_PRECISION (注3)	
MPI_2INTEGER	MPI_MAXLOC(最大と位置)
MPI_2REAL	MPI_MINLOC(最小と位置)
MPI_2DOUBLE_PRECISION(注5)	

MPIの基本データ型	MPIの定義済み演算
MPI_LOGICAL	MPI_LAND(論理AND)
	MPI_LOR (論理OR)
	MPI_LXOR(論理XOR)
MPI_INTEGER(注1)	MPI_BAND(ビットAND)
MPI_BYTE	MPI_BOR (ビットOR)
	MPI_BXOR(ビットXOR)

下記が使用できるマシンの環境もあります。

- (注1) MPI\_INTEGER4
- (注2) MPI\_REAL4
- (注3) MPI\_REAL8
- (注4) MPI\_COMPLEX16, MPI\_DOUBLE\_COMPLEX
- (注5) MPI\_2REAL8

### 使用法

```
CALL MPI_EXSCAN(sendbuf, recvbuf, count, datatype, op, comm, ierror)
```

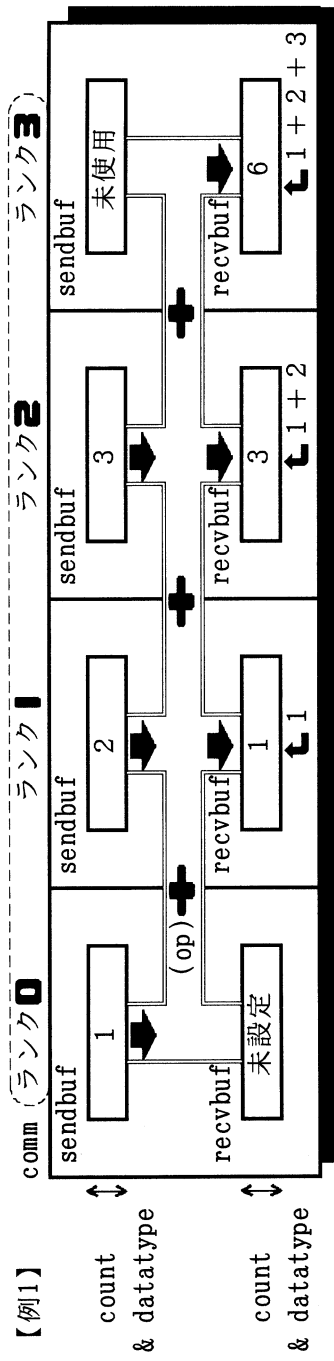
- sendbuf : 送信バッファァーの先頭アドレスを指定します。
- recvbuf : 受信バッファァーの先頭アドレスを指定します。  
送信バッファァーと受信バッファァーの実際に使用する部分は、メモリー上で重なってはいけません((MPI-2のMPI\_IN\_PLACEを使用することはできません))。
- count : 整数。送(受)メッセージの要素数を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- datatype : 整数。送(受)メッセージのデータ型を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- op : 整数。演算の種類を指定します。  
comm内の全プロセスが同じ値を指定する必要があります。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケータを指定します。  
全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

### 注意

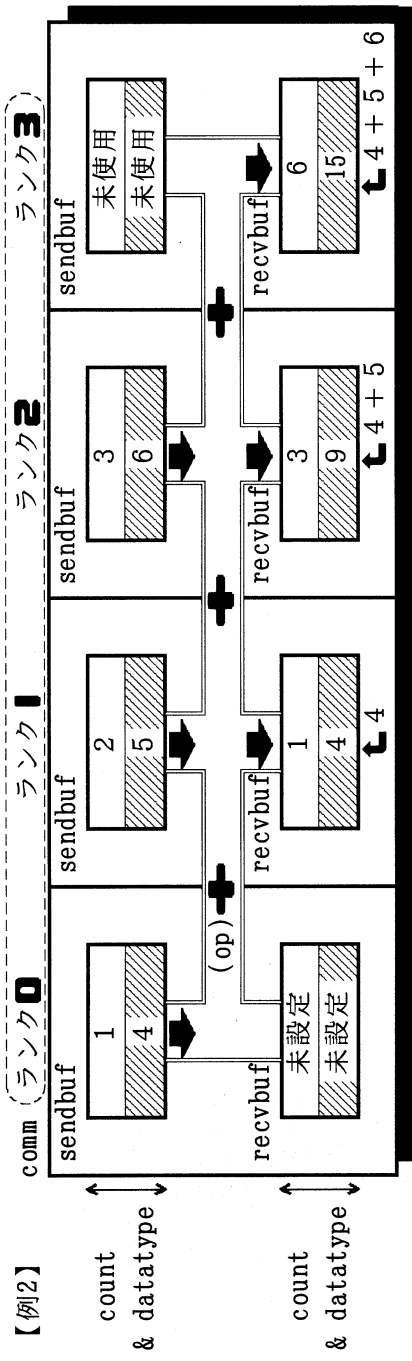
- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。



【例1】



【例2】



■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。
  ISEND = MYRANK + 1
  CALL MPI_EXSCAN(ISEND, IRECV, 1, MPI_INTEGER,
    & MPI_SUM, MPI_COMM_WORLD, IERR) 通信を行います。
  PRINT *, 'MYRANK = ', MYRANK, 'IRECV = ', IRECV 通信後の受信バッファの内容を書き出します。
  CALL MPI_FINALIZE(IERR)
  END

[ aoyama@node01 ] /u/aoyama: mpirun -np 4 a.out
MYRANK = 0 IRECV = (未設定)
MYRANK = 1 IRECV = 1
MYRANK = 2 IRECV = 3
MYRANK = 3 IRECV = 6
  
```

機能 (関連する節: 4-6-6-2節, 5-2-2節)

- コミュニケータ (comm) 内の、全プロセスの送信バッファ (sendbuf) のメッセージが、通信しながら演算 (op) されます。計算結果は、配列 recvcounts に指定した要素ずつ、ランクの小さい順に各プロセスの受信バッファ (recvbuf) に送信されます。例えば計算結果の最初から recvcounts(1) 個の要素はランク 0 のプロセスに、次の recvcounts(2) 個の要素はランク 1 のプロセスに...、のように送信されます。
- 送信バッファが配列の場合は、配列の対応する要素ごとに演算が行われます。
- 本サブルーチンは、メッセージを MPI\_REDUCE で 1 つのプロセスに収集した後、それを MPI\_SCATTERV で全プロセスに送信したのと機能的に類似しています。
- MPI の基本データ型 (datatype) と MPI の定義済み演算 (op) の組み合わせを以下に示します。演算 (op) をユーザーが自分で定義することもできます。

MPI の基本データ型	MPI の定義済み演算
MPI_INTEGER (注1)	MPI_SUM (合計)
MPI_REAL (注2)	MPI_PROD (積)
MPI_DOUBLE_PRECISION (注3)	
MPI_COMPLEX (注4)	
MPI_INTEGER (注1)	MPI_MAX (最大)
MPI_REAL (注2)	MPI_MIN (最小)
MPI_DOUBLE_PRECISION (注3)	
MPI_2INTEGER	MPI_MAXLOC (最大と位置)
MPI_2REAL	MPI_MINLOC (最小と位置)
MPI_2DOUBLE_PRECISION (注5)	

MPI の基本データ型	MPI の定義済み演算
MPI_LOGICAL	MPI_LAND (論理AND)
	MPI_LOR (論理OR)
	MPI_LXOR (論理XOR)
MPI_INTEGER (注1)	MPI_BAND (ビットAND)
MPI_BYTE	MPI_BOR (ビットOR)
	MPI_BXOR (ビットXOR)

下記が使用できるマシン環境もあります。

- (注1) MPI\_INTEGER4
- (注2) MPI\_REAL4
- (注3) MPI\_REAL8
- (注4) MPI\_COMPLEX16, MPI\_DOUBLE\_COMPLEX
- (注5) MPI\_2REAL8

使用法

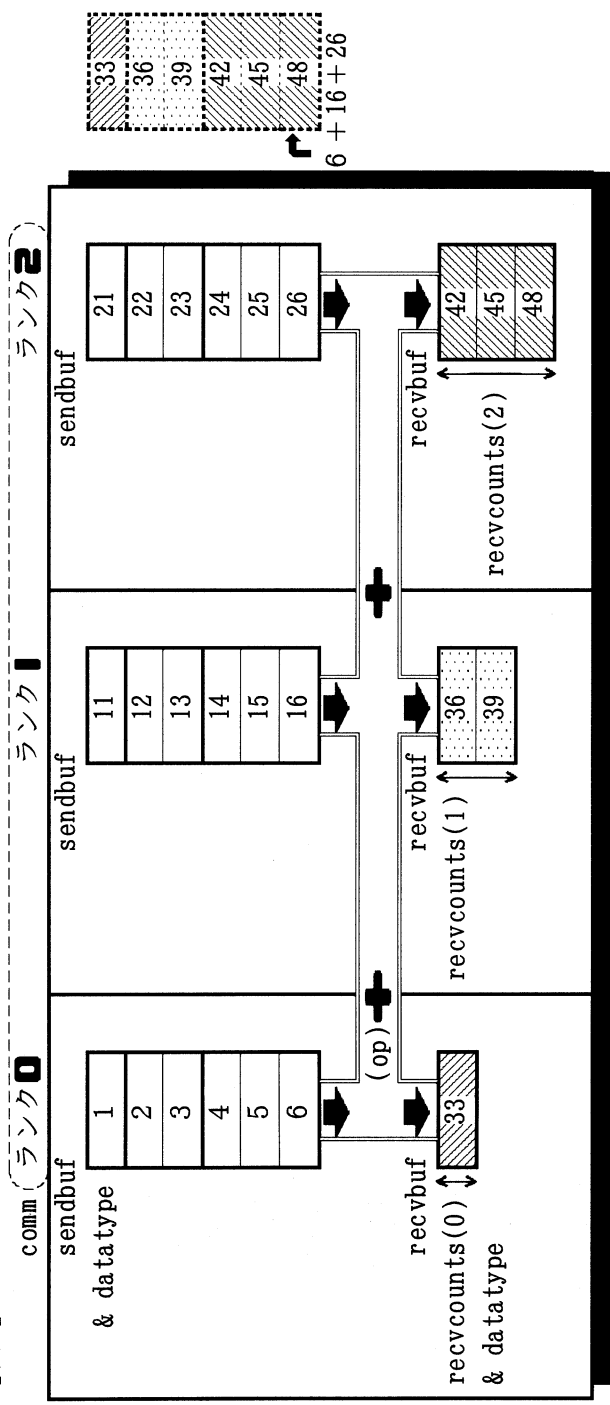
```
CALL MPI_REDUCE_SCATTER(sendbuf, recvbuf, recvcounts, datatype, op, comm, ierror)
```

- sendbuf : 送信バッファの先頭アドレスを指定します。
- recvbuf : 受信バッファの先頭アドレスを指定します。  
送信バッファと受信バッファの実際に使用する部分は、メモリー上で重なってはいけません(ただし MPI-2 の新機能 MPI\_IN\_PLACE で回避可能です。詳細は 3-8-1 節を参照)。
- recvcounts : 整数配列。ランク i のプロセスが受信するメッセージの要素数を、この配列の i+1 番目の要素に指定します。comm 内の全プロセスが同じ値を指定する必要があります(サンプルプログラムに示すように、配列を 0 から開始した方が分かりやすくなります)。
- datatype : 整数。送(受)信メッセージのデータ型を指定します。  
comm 内の全プロセスが同じ値を指定する必要があります。
- op : 整数。演算の種類を指定します。  
comm 内の全プロセスが同じ値を指定する必要があります。
- comm : 整数。送受信に参加する全てのプロセスを含むグループのコミュニケーションを指定します。  
全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

注意

- コミュニケータ (comm) 内の全プロセスが本サブルーチンをコールする必要があります。

【例1】



■ サンプルプログラム 【例1】のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISEND(6), IRECV(3)
  INTEGER IRCNT(0:2)
  DATA IRCNT/1,2,3/
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)  自分のランクMYRANKを取得します。
  DO I = 1, 6
    ISEND(I) = I + MYRANK*10
  ENDDO
  CALL MPI_REDUCE_SCATTER(ISEND,IRECV,IRCNT,MPI_INTEGER,
    & MPI_SUM,MPI_COMM_WORLD,IERR)
  PRINT *, 'MYRANK = ',MYRANK,' IRECV = ',IRECV  通信後の受信バッファの内容を書き出します。
  CALL MPI_FINALIZE(IERR)
  END

[aoyama@node01]/u/aoyama: mpirun -np 3 a.out
MYRANK = 0 IRECV = 33 0 0
MYRANK = 1 IRECV = 36 39 0
MYRANK = 2 IRECV = 42 45 48
  
```

# 集団通信サブルーチン MPI\_OP\_CREATE

## 機能 (関連する節: 3-3-6-2節)

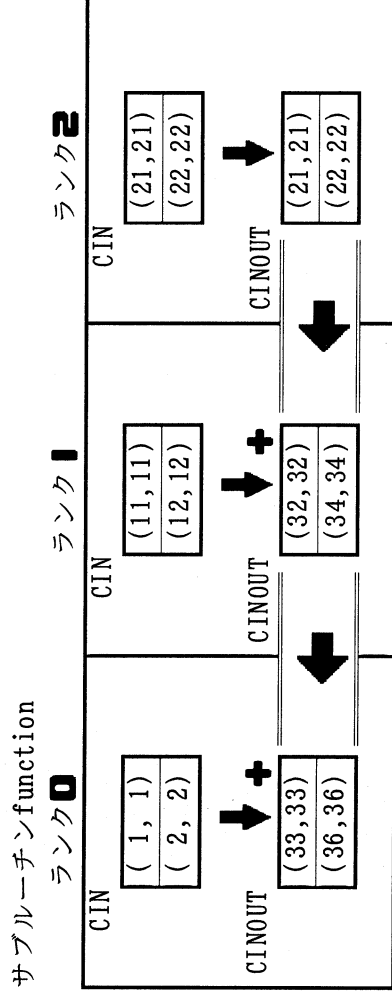
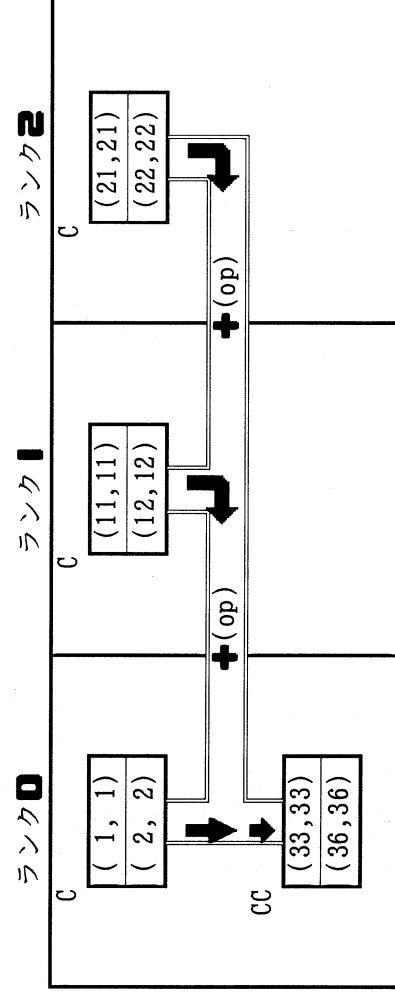
MPI\_REDUCE(V), MPI\_ALLREDUCE(V), MPI\_SCAN(V), MPI\_REDUCE\_SCATTER(V)で使用する演算(op)をユーザがサブルーチンfunctionに作成した場合、本サブルーチンはその演算を登録し、演算のIDがopに戻ります。本サブルーチンを実行した後、opをMPI\_REDUCEなどの引数opとして使用することができます。なお、ユーザが作成するサブルーチン(function)の指定方法については3-3-6-2節を参照して下さい。

## 使用法

CALL MPI\_OP\_CREATE(function, commute, op, ierror)

- function: ユーザが演算を定義したサブルーチン名を指定します。MPI\_OP\_CREATEをコールしたサブルーチンでは、ここで指定したサブルーチン名を外部関数としてEXTERNAL文で宣言して下さい。
- commute: ユーザが定義した演算は結合則を満たす必要がありますが、それに加えて交換則を満たす場合は.TRUE.を、満たさない場合は.FALSE.を指定して下さい。.TRUE.を指定した場合は演算(通信)順序が状況に応じて最適化され、.FALSE.を指定した場合にはプロセスのランクが0から昇順に演算(通信)が行われます(3-3-6節の【例1】参照)。
- op: 整数。登録された演算のIDが戻ります。この値を保存しておき、MPI\_REDUCEなどの引数opで使用します。MPI\_OP\_CREATEをコールしたサブルーチン以外のサブルーチンでMPI\_REDUCEなどをコールするときは、この値をCOMMON文などで渡す必要があります。
- ierror: 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

【例1】倍精度複素数の加算の例です(3-3-6-2節参照)。



## ■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  EXTERNAL PARA_SUM16
  COMPLEX*16 C(2),CC(2)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)  プロセス数NPROCSを取得します。
  CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)  自分のランクMYRANKを取得します。
  CALL MPI_OP_CREATE(PARA_SUM16,.TRUE.,ISUM16,IERR) 演算を登録し、演算のIDがISUM16に戻ります。
  C(1) = (MYRANK*10+1,MYRANK*10+1) 各プロセスは送信バッファアーにデータをセットします。
  C(2) = (MYRANK*10+2,MYRANK*10+2)
  CALL MPI_REDUCE(C,CC,2,MPI_DOUBLE_COMPLEX,ISUM16,0, ISUM16)を使用して通信を行います。
  & MPI_COMM_WORLD, IERR)
  IF (MYRANK == 0) PRINT *, 'C =', CC 通信後の受信バッファアーの内容を書き出します。
  CALL MPI_FINALIZE(IERR)  MPI環境の終了処理を行います。
  END

SUBROUTINE PARA_SUM16(CIN,CINOUT,LEN,ITYPE) ユーザーが定義した演算(倍精度複素数の加算)です。
  COMPLEX*16 CIN(*),CINOUT(*)
  DO I = 1, LEN
    CINOUT(I) = CIN(I) + CINOUT(I)
  ENDDO
  END

```

```

[aoyama@node01]~/u/aoyama: mpirun -np 3 a.out
C = (33.000...,33.000...) (36.000...,36.000...)

```

# 集団通信サブルーチン MPI\_BARRIER

## 機能 (関連する節: 4-8-3節、4-4-2節)

コミュニケータ(comm)内の全プロセス間で(何らかの理由で)同期をとりたい場合、本サブルーチンを使用します。

図1のように『処理1』➡『CALL MPI\_BARRIER』➡『処理2』の順に処理を行うプログラムを3プロセスで並列に実行する場合を想定します。このとき『処理1』の部分の進行が、図2のようにプロセスによって異なっているとしています。

まず、ランク0のプロセスが『CALL MPI\_BARRIER』に到達しますが、他のプロセスがまだ到達していませんため、ランク0はウェイトします。次にランク1が『CALL MPI\_BARRIER』に到達し、同様にウェイトします。そして最後にランク2が『CALL MPI\_BARRIER』に到達すると、各プロセスはいっせいに実行を再開します。

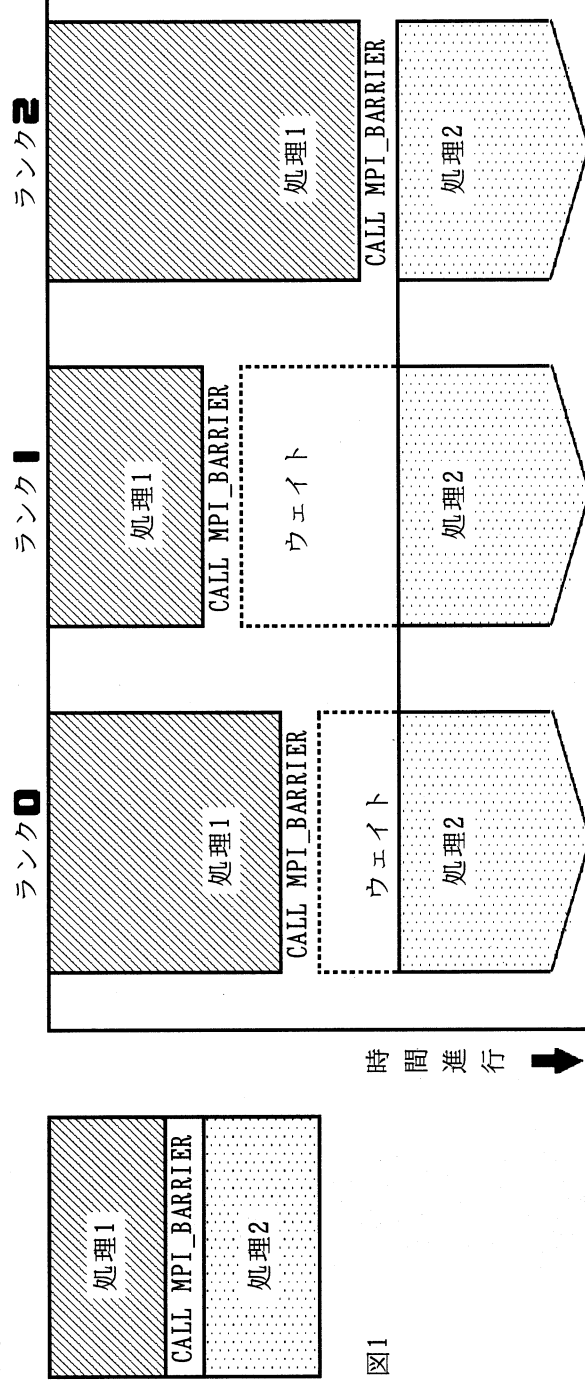


図2

## 使用法

```
CALL MPI_BARRIER(comm, ierror)
```

- comm : 整数。同期をとりたい全てのプロセスを含むグループのコミュニケータを指定します。全プロセスが同じ値を指定する必要があります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## 注意

- コミュニケータ(comm)内の全プロセスが本サブルーチンをコールする必要があります。

## ■ サンプルプログラム

『CALL MPI\_BARRIER』を、並列ジョブの経過時間測定の誤差を少なくするために用いる例を以下に示します。

以下の例で、『計算部分』の経過時間を測定するとします。ところが、『入力データの読み込みなど』の部分は、各プロセスで経過時間にばらつきがあるため(特に共有ディスクから読み込みを行う場合)、測定の開始時点(②)がプロセスによって異なってしまいます。そこで①に『CALL MPI\_BARRIER』を挿入し、プロセス間の同期をとります。

『計算部分』の終了後、④で測定を終了します。通常『計算部分』の経過時間はプロセス間でばらつきがありますが、最も長くかかったプロセスの経過時間を知りたい場合は③にも『CALL MPI\_BARRIER』を挿入してプロセス間の同期をとります。一方、各プロセスの『計算部分』の経過時間のばらつきを比較したい場合は③を指定しません。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  REAL*8 ELP1, ELP2
  CALL MPI_INIT(IERR)

  入力データの読み込みなど

  CALL MPI_BARRIER(MPI_COMM_WORLD,IERR) ① プロセス間の同期をとります。
  ELP1 = MPI_WTIME()                       ② MPI_WTIME()を実行して測定を開始します。

  計算部分

  CALL MPI_BARRIER(MPI_COMM_WORLD,IERR) ③ プロセス間の同期をとります。
  ELP2 = MPI_WTIME()                       ④ MPI_WTIME()を実行して測定を終了します。
  PRINT *, 'ELAPSE = ', ELP2-ELP1
  CALL MPI_FINALIZE(IERR)
  END
  
```

MPIが使用するインクルードファイルを指定します。  
経過時間の測定に使用する変数を倍精度で宣言します。  
MPI環境の初期化処理を行います。

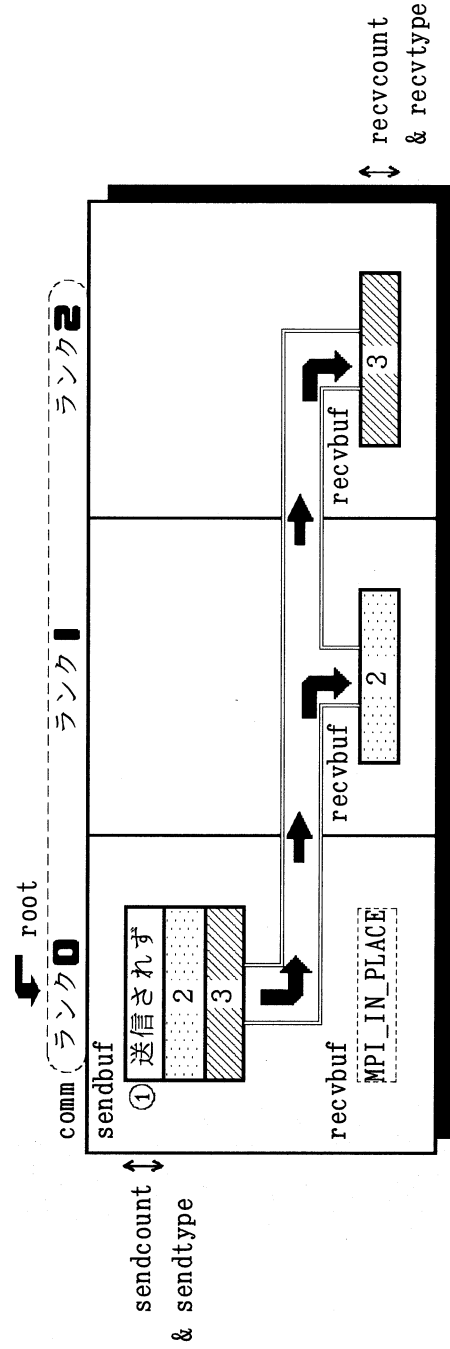
## 【MPI-2】MPI\_IN\_PLACEの使い方 MPI\_SCATTER

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプションをサポートしていない場合は、使用することができません(3-8-1節参照)。

### ■ 使用法

```
CALL MPI_SCATTER(sendbuf, sendcount, sendtype,
                recvbuf, recvcnt, recvtype, root, comm, ierror)
```

- 送信元プロセス(root)では、recvbufにMPI\_IN\_PLACEを指定します。rootプロセスではrecvcountとrecvtypeは無視されますので、適当な変数を指定してもかまいません。
- 下図の①の部分(MPI\_IN\_PLACEを使用しない場合にrootプロセスの送信データを入れる部分)は送信されません。



■ サンプルプログラム 上図のプログラム例です。

```

:
INTEGER ISEND(3)
:
IF (MYRANK==0) THEN
CALL MPI_SCATTER(ISEND ,1 ,MPI_INTEGER,
                MPI_IN_PLACE, IDUMMY, IDUMMY ,
                0, MPI_COMM_WORLD, IERR)
&
&
ELSE
CALL MPI_SCATTER(IDUMMY, IDUMMY, IDUMMY ,
                IRECV ,1 ,MPI_INTEGER,
                0, MPI_COMM_WORLD, IERR)
&
&
ENDIF
:
rootプロセス
その他のプロセス
```

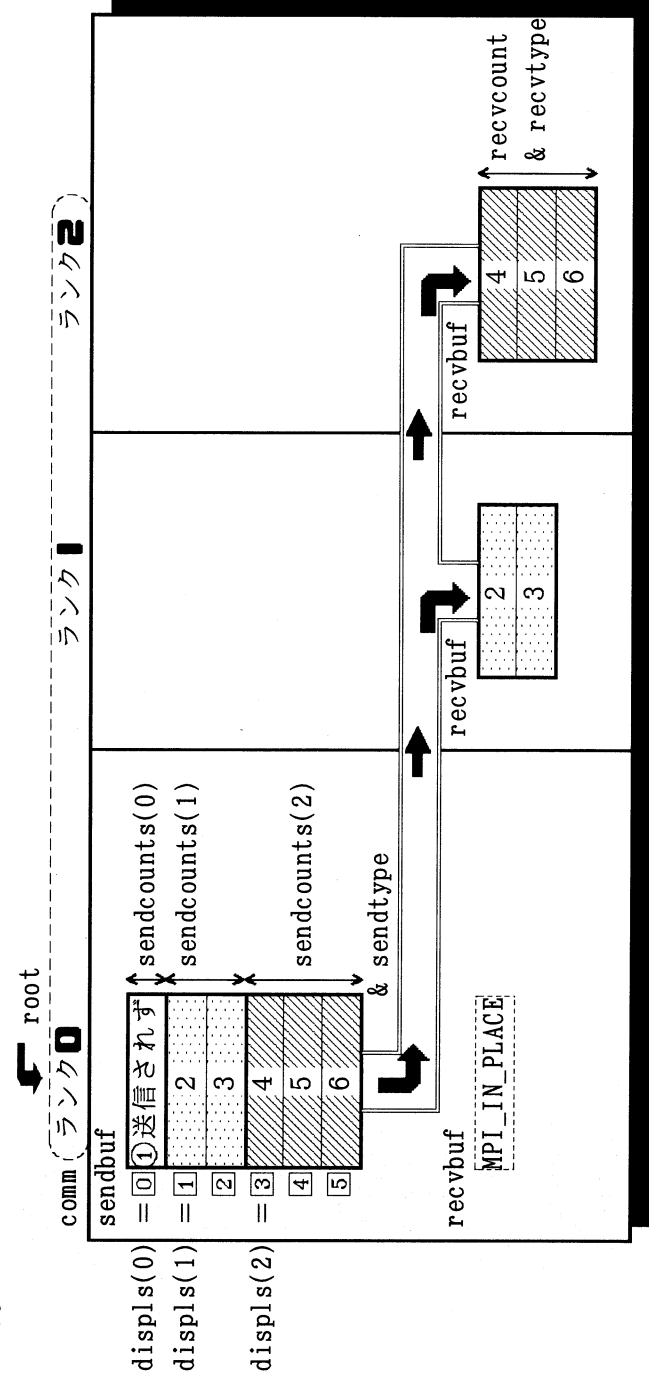


以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。  
この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプシヨンをサポートしていない場合は、使用することができません(3-8-1節参照)。

■ 使用法

```
CALL MPI_SCATTERV(sendbuf, sendcounts, displs, sendtype,
                 recvbuf, recvcount, recvtype, root, comm, ierror)
```

- 送信元プロセス(root)では、recvbufにMPI\_IN\_PLACEを指定します。rootプロセスではrecvcountとrecvtypeは無視されますので、適当な変数を指定してもかまいません。
- 下図の①の部分(MPI\_IN\_PLACEを使用しない場合にrootプロセスの送信データを入れる部分)は送信されません。



■ サンプルプログラム 上図のプログラム例です。

```
INTEGER ISEND(6), IRECV(3), ISCNT(0:2), IDISP(0:2)
DATA ISCNT/1,2,3/ IDISP/0,1,3/
IF (MYRANK==0) THEN
CALL MPI_SCATTERV( ISEND, ISCNT, IDISP, MPI_INTEGER,
                  MPI_IN_PLACE, IDUMMY, IDUMMY, IDUMMY, IDUMMY, IDUMMY,
                  0, MPI_COMM_WORLD, IERR )
ELSE
CALL MPI_SCATTERV( IDUMMY, IDUMMY, IDUMMY, IDUMMY, IDUMMY, IDUMMY,
                  IRECV, MYRANK+1, MPI_INTEGER,
                  0, MPI_COMM_WORLD, IERR )
ENDIF
```

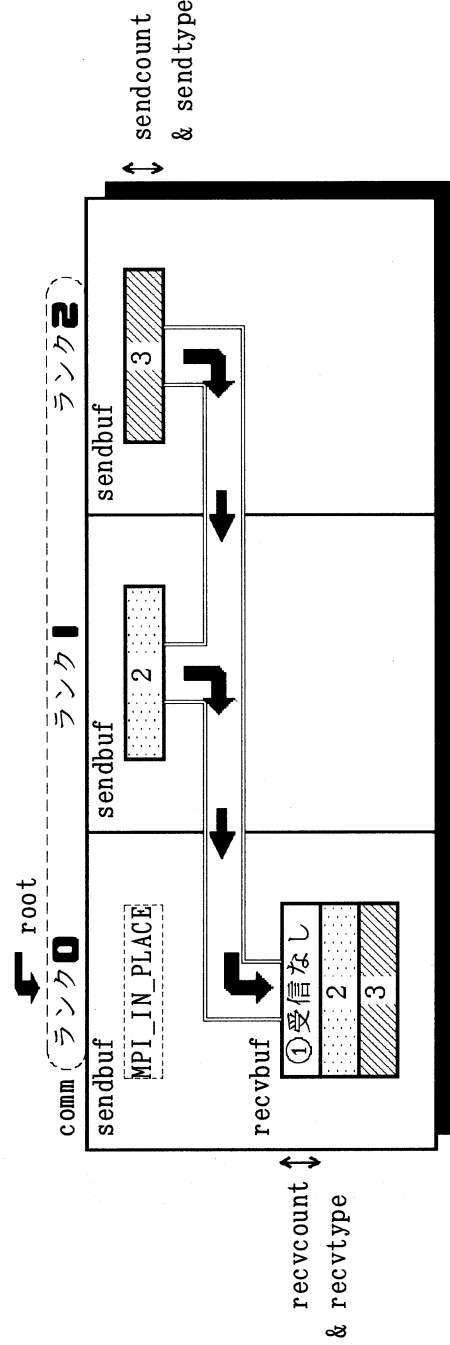
## 【MPI-2】MPI\_IN\_PLACEの使い方 MPI\_GATHER

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプションをサポートしていない場合は、使用することができません(3-8-1節参照)。

### ■ 使用法

```
CALL MPI_GATHER(sendbuf, sendcount, sendtype,
               recvbuf, recvcount, recvtype, root, comm, ierror)
```

- 宛先プロセス(root)では、sendbufにMPI\_IN\_PLACEを指定します。rootプロセスではsendcountとsendtypeは無視されますので、適当な変数を指定してもかまいません。
- 下図の①の部分(MPI\_IN\_PLACEを使用しない場合にrootプロセスの受信データが入る部分)には受信されません。



■ サンプルプログラム 上図のプログラム例です。

```

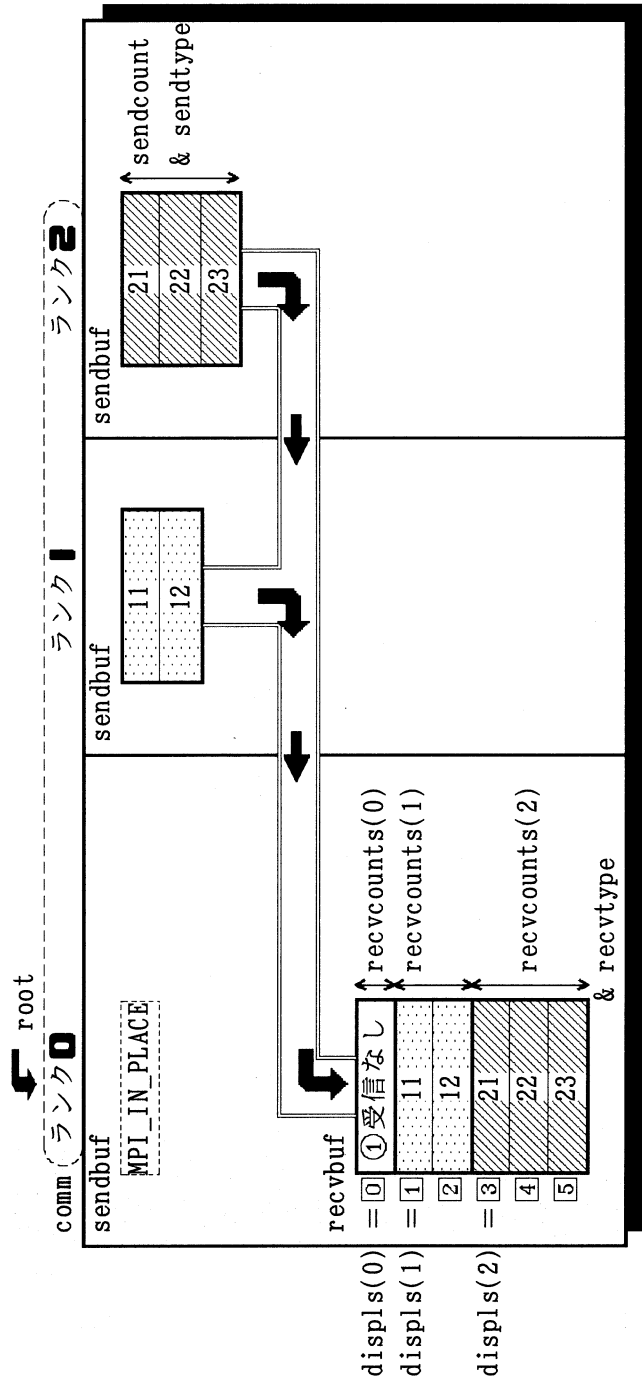
:
INTEGER IRECV(3)
:
IF (MYRANK==0) THEN
  CALL MPI_GATHER(MPI_IN_PLACE, IDUMMY, IDUMMY, IDUMMY, IDUMMY,
                 IRECV, 1, MPI_INTEGER,
                 0, MPI_COMM_WORLD, IERR)
ELSE
  CALL MPI_GATHER(IRECV, 1, MPI_INTEGER,
                 IDUMMY, IDUMMY, IDUMMY, IDUMMY,
                 0, MPI_COMM_WORLD, IERR)
ENDIF
:
rootプロセス
その他のプロセス
```

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプシヨンをサポートしていない場合は、使用することができません(3-8-1節参照)。(関連する節: 4-6-3-1節)

### ■ 使用法

```
CALL MPI_GATHERV(sendbuf, sendcount, sendtype,
                 recvbuf, recvcounts, displs, recvtype, root, comm, ierror)
```

- 宛先プロセス (root) では、sendbuf に MPI\_IN\_PLACE を指定します。root プロセス では sendcount と sendtype は無視されますので、適当な変数を指定してもかまいません。
- 下図の①の部分 (MPI\_IN\_PLACE を使用しない場合に root プロセス の受信データが入る部分) には受信されません。



■ サンプルプログラム 上図のプログラム例です。

```

:
INTEGER ISEND(3), IRECV(6), IRCNT(0:2), IDISP(0:2)
DATA IRCNT/1,2,3/ IDISP/0,1,3/
:
IF (MYRANK==0) THEN
  CALL MPI_GATHERV(MPI_IN_PLACE, IDUMMY, IDUMMY,
                  IRECV, IRCNT, IDISP, MPI_INTEGER,
                  0, MPI_COMM_WORLD, IERR)
  rootプロセス
ELSE
  CALL MPI_GATHERV(ISEND, MYRANK+1, MPI_INTEGER,
                  IDUMMY, IDUMMY, IDUMMY, IDUMMY,
                  0, MPI_COMM_WORLD, IERR)
  その他のプロセス
ENDIF
:

```

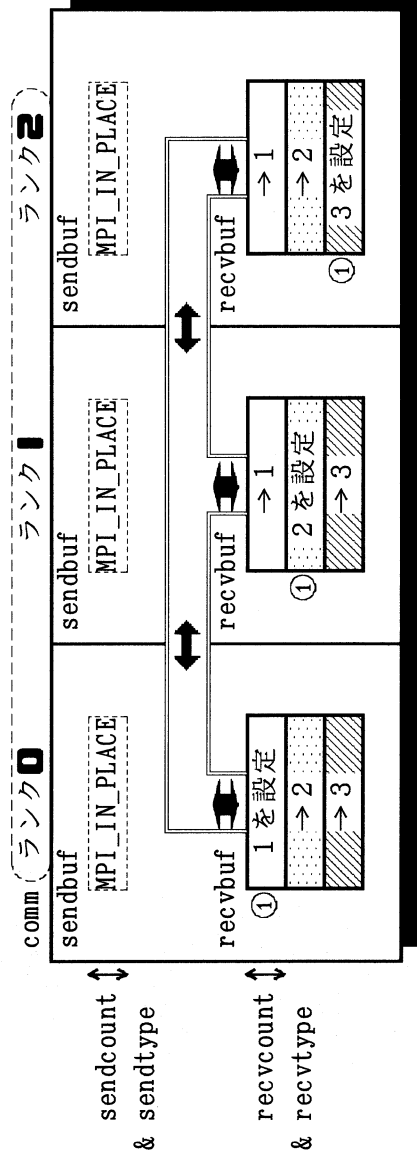
## [MPI-2] MPI\_IN\_PLACEの使い方 MPI\_ALLGATHER

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプションをサポートしていない場合は、使用することができません(3-8-1節参照)。

### ■ 使用法

```
CALL MPI_ALLGATHER(sendbuf, sendcount, sendtype,
                  recvbuf, recvcount, recvtype, comm, ierror)
```

- 全プロセスで、sendbufにMPI\_IN\_PLACEを指定します。sendcountとsendtypeは無視されますので、適当な変数を指定してもかまいません。
- 各プロセスは下図の①に送信データを設定します。通信終了後、下図の→のように変化します。



■ サンプルプログラム 上図のプログラム例です。

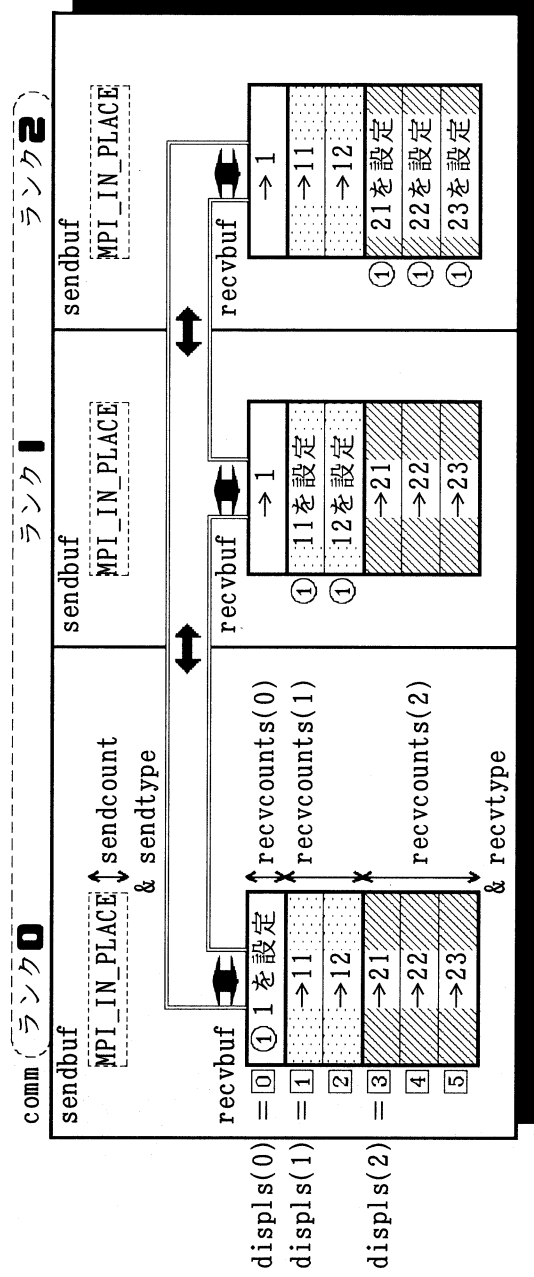
```
INTEGER IRECV(3)
:
CALL MPI_ALLGATHER(MPI_IN_PLACE, IDUMMY, IDUMMY,
                  IRECV, 1, MPI_INTEGER,
                  MPI_COMM_WORLD, IERR)
:
```

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。  
この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプシヨンをサポートしていない場合は、使用することができません(3-8-1節参照)。(関連する節:4-6-3-2節)

### ■ 使用法

```
CALL MPI_ALLGATHERV(sendbuf, sendcount, sendtype,
                    recvbuf, recvcounts, displs, recvtype, comm, ierror)
```

- 全プロセスで、sendbufにMPI\_IN\_PLACEを指定します。sendcountとsendtypeは無視されますので、適当な変数を指定してもかまいません。
- 各プロセスは下図の①に送信データを設定します。通信終了後、下図の②のように変化します。



■ サンプルプログラム 上図のプログラム例です。

```

:
INTEGER IRECV(6), IRCNT(0:2), IDISP(0:2)
DATA IRCNT/1,2,3/ IDISP/0,1,3/
:
CALL MPI_ALLGATHERV(MPI_IN_PLACE, IDUMMY, IDUMMY,
& IRECV, IRCNT, IDISP, MPI_INTEGER,
& MPI_COMM_WORLD, IERR)
:

```

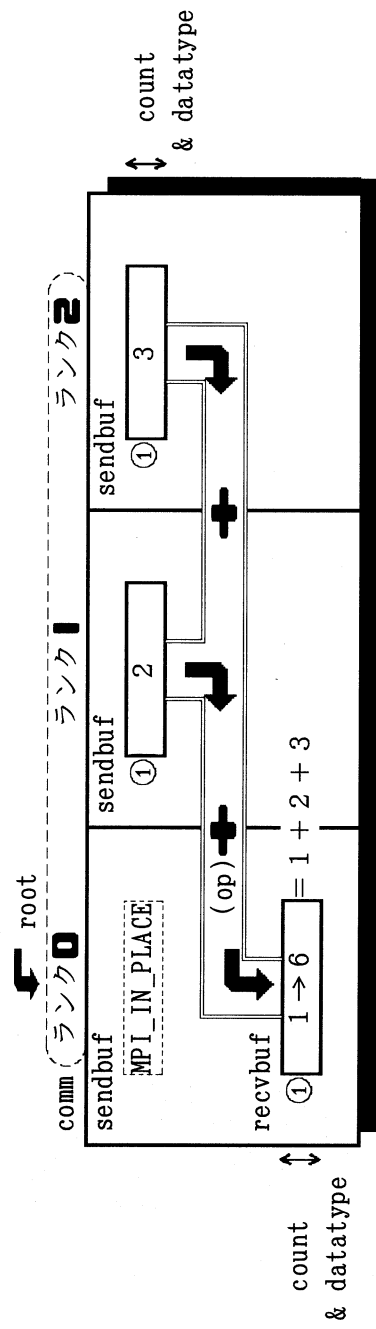
## [MPI-2] MPI\_IN\_PLACEの使い方 MPI\_REDUCE

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプションをサポートしていない場合は、使用することができません(3-8-1節参照)。

### ■ 使用法

```
CALL MPI_REDUCE(sendbuf, recvbuf, count, datatype, op, root, comm, ierror)
```

- 宛先プロセス(root)では、sendbufにMPI\_IN\_PLACEを指定します。
- 各プロセスは下図の①に送信データを設定します(ランク0の①の位置に注意して下さい)。通信終了後、下図の→のように変化します。



■ サンプルプログラム 上図のプログラム例です。

```

:
IF (MYRANK==0) THEN
CALL MPI_REDUCE(MPI_IN_PLACE, IRECV, 1, MPI_INTEGER,
& MPI_SUM, 0, MPI_COMM_WORLD, IERR)
ELSE
CALL MPI_REDUCE(ISEND, IDUMMY, 1, MPI_INTEGER,
& MPI_SUM, 0, MPI_COMM_WORLD, IERR)
ENDIF
:
その他のプロセス

```

## 【MPI-2】MPI\_IN\_PLACEの使い方 MPI\_ALLREDUCE

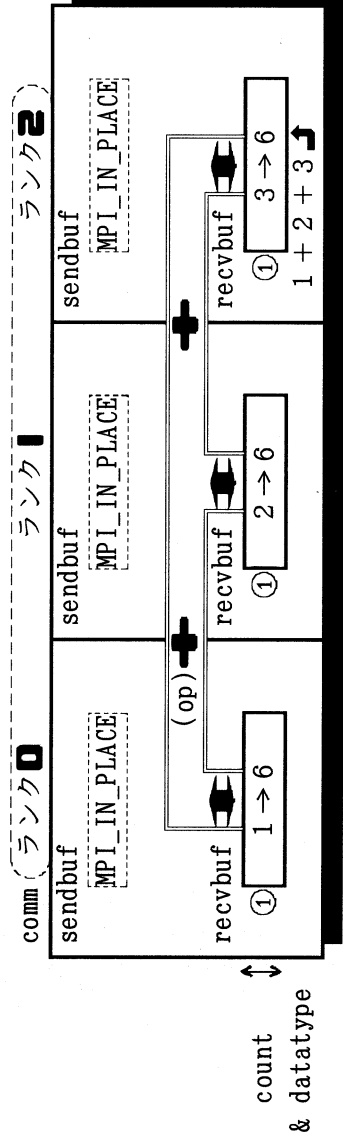
付録 MPI サブルーチン一覧

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプションをサポートしていない場合は、使用することができません(3-8-1節参照)。

### ■ 使用法

```
CALL MPI_ALLREDUCE(sendbuf, recvbuf, count, datatype, op, comm, ierror)
```

- 全プロセスで、sendbufにMPI\_IN\_PLACEを指定します。
- 各プロセスは下図の①に送信データを設定した後、下図の→のように変化します。



■ サンプルプログラム 上図のプログラム例です。

```

:
CALL MPI_ALLREDUCE(MPI_IN_PLACE, IRECV, 1, MPI_INTEGER,
& MPI_SUM, MPI_COMM_WORLD, IERR)
:

```

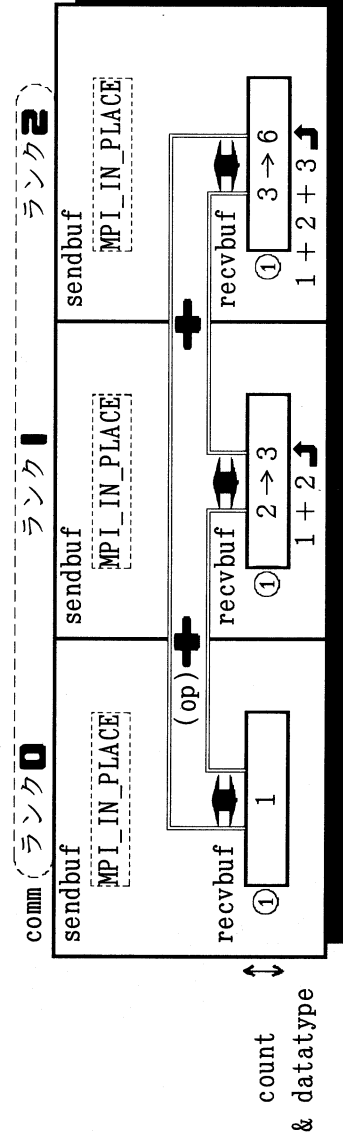
## 【MPI-2】MPI\_IN\_PLACEの使い方 MPI\_SCAN

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプションをサポートしていない場合は、使用することができません(3-8-1節参照)。

### ■ 使用法

```
CALL MPI_SCAN(sendbuf, recvbuf, count, datatype, op, comm, ierror)
```

- 全プロセスで、sendbufにMPI\_IN\_PLACEを指定します。
- 各プロセスは下図の①に送信データを設定します。通信終了後、下図の→のように変化します。



■ サンプルプログラム 上図のプログラム例です。

```
CALL MPI_SCAN(MPI_IN_PLACE, IRECV, 1, MPI_INTEGER,  
& MPI_SUM, MPI_COMM_WORLD, IERR)
```



## 【MPI-2】MPI\_IN\_PLACEの使い方 MPI\_REDUCE\_SCATTER

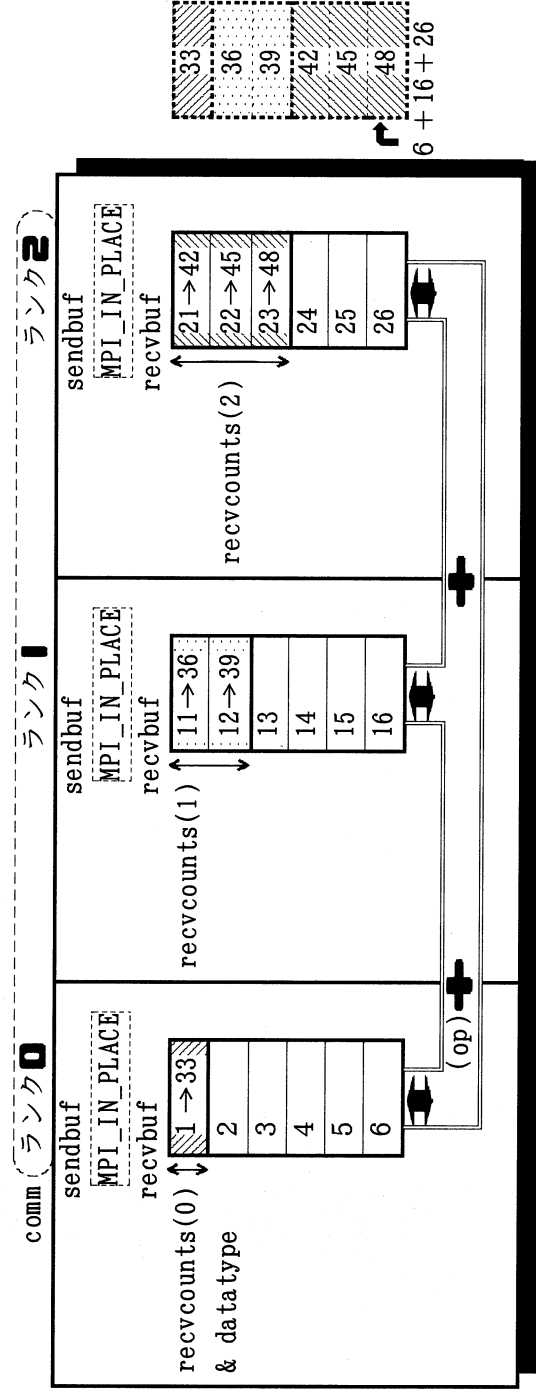
付録 MPI サブルーチン一覽

以下では、当MPIサブルーチンの機能のうち、MPI\_IN\_PLACEに関連する部分のみを説明します。  
この機能は、使用するマシン環境がMPI-2のMPI\_IN\_PLACEオプションをサポートしていない場合は、使用することができません(3-8-1節参照)。

### ■ 使用法

```
CALL MPI_REDUCE_SCATTER(sendbuf, recvbuf, recvcnts, datatype, op, comm, ierror)
```

- 全プロセスで、sendbufにMPI\_IN\_PLACEを指定します。
- 各プロセスは下図のように送信データを設定します。通信終了後、送信されたデータは(全プロセスで)受信バッファrecvbufの先頭から入り、図の→のように変化します。



■ サンプルプログラム 上図のプログラム例です。

```

:
INTEGER IRECV(6), IRCNT(0:2)
DATA IRCNT/1, 2, 3/
:
CALL MPI_REDUCE_SCATTER(MPI_IN_PLACE, IRECV, IRCNT, MPI_INTEGER,
& MPI_SUM, MPI_COMM_WORLD, IERR)
:

```

# 1対1ブロッキング通信サブルーチン MPI\_SEND

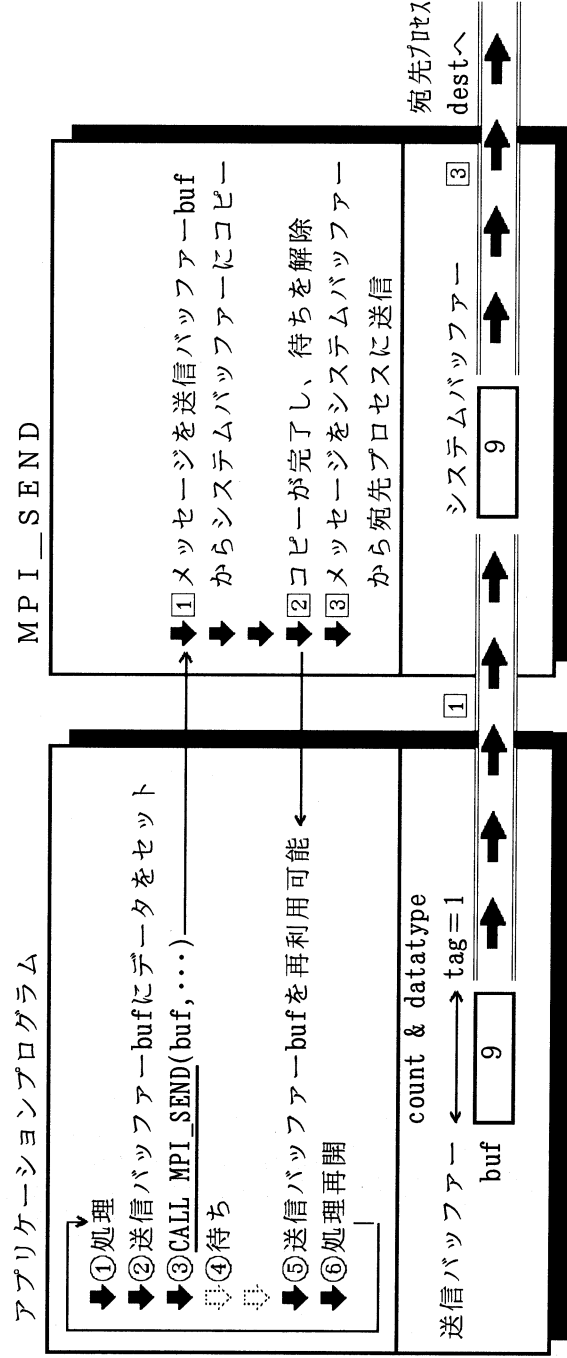
## 機能 (関連する節: 3 - 4 節)

送信バッファ (buf) 内の、データ型が `datatype` である `count` 個の連続した要素から構成した送信メッセージを、タグ (`tag`) をつけて宛先プロセス (コミュニケーションでランクが `dest`) に送信します。MPI\_SEND で送信したメッセージは、MPI\_RECV と MPI\_IRECV のどちらかで受信してもかまいません。

MPI\_SEND をコールした後、アプリケーションプログラムは、送信バッファが再利用可能になるまで待ちの状態に入ります。

## 動作

- アプリケーションプログラムは、②で送信バッファ `buf` にデータ (メッセージ) をセットし、③で MPI\_SEND をコールした後、④で待ちの状態に入ります。
- サブルーチン MPI\_SEND は、①でメッセージを送信バッファ `buf` からシステムバッファにコピーします。コピーが完了した後、②でアプリケーションプログラムの待ちを解除します。なお、待ちが解除されたということは、送信したメッセージが宛先プロセスに到着したことを意味するわけではありません。
- ③の時点で、宛先プロセスが MPI\_RECV または MPI\_IRECV をすでにコールしている場合は、システムバッファのメッセージを宛先プロセスに送信します。宛先プロセスがまだコールしていない場合は、コールするまで③で待ちます。
- アプリケーションプログラムは、⑤で送信バッファ `buf` の再利用が可能になり、⑥で処理を再開します。
- ②で再び送信バッファ `buf` にデータ (メッセージ) をセットしても、前回の②で送信バッファ `buf` にセットしたメッセージはすでに②でシステムバッファにコピーされているので、壊れることはありません。



## ■ 使用法

```
CALL MPI_SEND(buf, count, datatype, dest, tag, comm, ierror)
```

- buf : 送信バッファの先頭アドレスを指定します。
- count : 整数。送信バッファ内の送信メッセージの要素数を指定します。
- datatype : 整数。送信バッファ内の送信メッセージのデータ型を指定します。
- dest : 整数。宛先プロセスのcomm内でのランクを指定します。宛先プロセスに対してメッセージを送信しない場合には、MPI\_PROC\_NULLを指定します。
- tag : 整数。送信するメッセージの種類を区別したい場合に使用します。特に区別する必要がない場合には、適当な値(例えば1)を指定します。
- 0 ~ MPI\_TAG\_UB までの整数を指定することができます。MPI\_TAG\_UBの値はmpi.f.hの中で設定されており、マシン環境によって異なります(PRINT文で書き出せば値が分かります)。
- comm : 整数。自分のプロセスと宛先プロセスを含むグループのコミュニケーションを指定します。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

■ サンプルプログラム 【例1】 のプログラム例です。ランク0からランク1にメッセージを送信します。

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)  プロセス数NPROCSを取得します。
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)  自分のランクMYRANKを取得します。
  IF (MYRANK == 0) ISBUF = 9
  ランク0のプロセスは送信バッファにデータを
  セットします。

  IF (MYRANK == 0) THEN
    CALL MPI_SEND(ISBUF, 1, MPI_INTEGER,
                  1, 1, MPI_COMM_WORLD, IERR)
    ランク0はランク1にメッセージを送信します。
  &
  ELSEIF (MYRANK == 1) THEN
    CALL MPI_RECV(IRBUF, 1, MPI_INTEGER,
                  0, 1, MPI_COMM_WORLD, ISTATUS, IERR)
    ランク1はランク0からメッセージを受信します。
  &
  ENDIF
  IF (MYRANK == 1) PRINT *, 'IRBUF = ', IRBUF  通信後の受信バッファの内容を書き出します。
  CALL MPI_FINALIZE(IERR)
  END
```

```
[aoyama@node011]/u/aoyama: mpi.run -np 2 a.out
IRBUF = 9
```

# 1対1ブロッキング通信サブルーチン MPI\_RECV

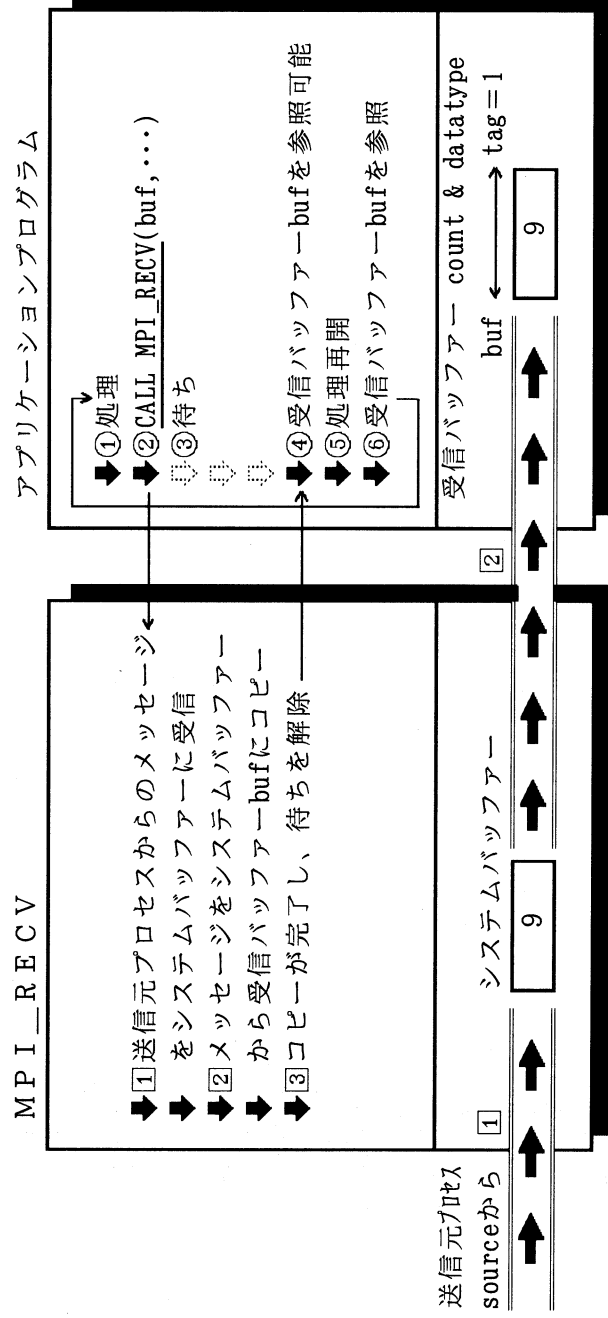
## 機能 (関連する節: 3 - 4 節)

送信元プロセス(コミュニケータcomm内でランクがsource)から送信された、タグがtagのメッセージを、データ型がdatatypeであるcount個の連続した要素から構成される受信バッファ( buf )に受信します。MPI\_SENDとMPI\_ISENDのどちらで送信したメッセージも、MPI\_RECVで受信することができます。

MPI\_RECVをコールした後、アプリケーションプログラムは、受信バッファへのメッセージの受信が完了するまで待ちの状態に入ります。

## 動作

- アプリケーションプログラムは、②でMPI\_RECVをコールした後、③で待ちの状態に入ります。
- サブルーチンMPI\_RECVは、①の時点で送信元プロセスがすでにMPI\_SENDまたはMPI\_ISENDをコールしている場合は、送信元プロセスから送られたメッセージをシステムバッファに受信します。送信元プロセスがまだコールしていない場合は、コールするまで④で待ちます。
- ②でメッセージをシステムバッファから受信バッファへコピーします。コピーが完了した後、③でアプリケーションプログラムの待ちを解除します。
- アプリケーションプログラムは、④で受信バッファ( buf )を参照可能になり、⑤で処理を再開します。
- ⑥で受信バッファ( buf )を参照したときには、受信メッセージが全て入っています。



## ■ 使用法

```
CALL MPI_RECV(buf, count, datatype, source, tag, comm, status, ierror)
```

- buf : 受信バッファの先頭アドレスを指定します。
- count : 整数。受信バッファの要素数を指定します。受信メッセージの要素数と同じか、それ以上の大きさにして下さい。受信メッセージの要素数より少ない場合は、エラーになります(マシン環境によってはエラーにならない場合もあります)。
- datatype : 整数。受信バッファに入る受信メッセージのデータ型を指定します。
- source : 整数。受信したいメッセージを送信するプロセスのcomm内でのランクを指定します。任意のプロセスからのメッセージを受信したい場合は、ワールドカードMPI\_ANY\_SOURCEを指定します。任意のプロセスからのメッセージを受信しない場合には、MPI\_PROC\_NULLを指定します。
- tag : 整数。受信したいメッセージに付けられているタグの値を指定します。この値は送信元プロセスがそのメッセージを送信した際に引数tagとして付けた値です。  
任意のタグ値のメッセージを受信したい場合は、ワールドカードMPI\_ANY\_TAGを指定します。
- comm : 整数。自分のプロセスと送信元プロセスを含むグループのコミュニケーションを指定します。
- status : 状況オブジェクトの配列を指定します。サンプルプログラムに示すように、大きさがMPI\_STATUS\_SIZEの整数配列を、『INCLUDE 'mpif.h'』より後に指定して下さい。本サブルーチンが完了すると、受信したメッセージの送信元ランクがstatus(MPI\_SOURCE)に、タグがstatus(MPI\_TAG)に戻ります。これらの値は通常は使用しません。使用するのは、ワールドカードMPI\_ANY\_SOURCEやMPI\_ANY\_TAGを指定した際、受信したメッセージの送信元ランクやタグの値を知りたい場合です(4-7-1-3節参照)。その場合、受信したメッセージの要素数も知る必要があるならば、『CALL MPI\_GET\_COUNT』を使用して下さい。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

■ サンプルプログラム 【例1】 のプログラム例です。ランク 1 はランク 0 からメッセージを受信します。

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  IF (MYRANK == 0) ISBUF = 9
  &
  IF (MYRANK == 0) THEN
    CALL MPI_SEND(ISBUF, 1, MPI_INTEGER,
                  1, 1, MPI_COMM_WORLD, IERR)
  &
  ELSEIF (MYRANK == 1) THEN
    CALL MPI_RECV(IRBUF, 1, MPI_INTEGER,
                  0, 1, MPI_COMM_WORLD, ISTATUS, IERR)
  &
  ENDIF
  IF (MYRANK == 1) PRINT *, IRBUF = ', IRBUF
  CALL MPI_FINALIZE(IERR)
END
```

MPIが使用するインクルードファイルを指定します。  
MPI\_RECVで使用する状況オブジェクトを指定します。  
MPI環境の初期化処理を行います。  
プロセス数NPROCSを取得します。  
自分のランクMYRANKを取得します。  
ランク0のプロセスは送信バッファにデータをセットします。  
ランク0はランク1にメッセージを送信します。  
ランク1はランク0からメッセージを受信します。

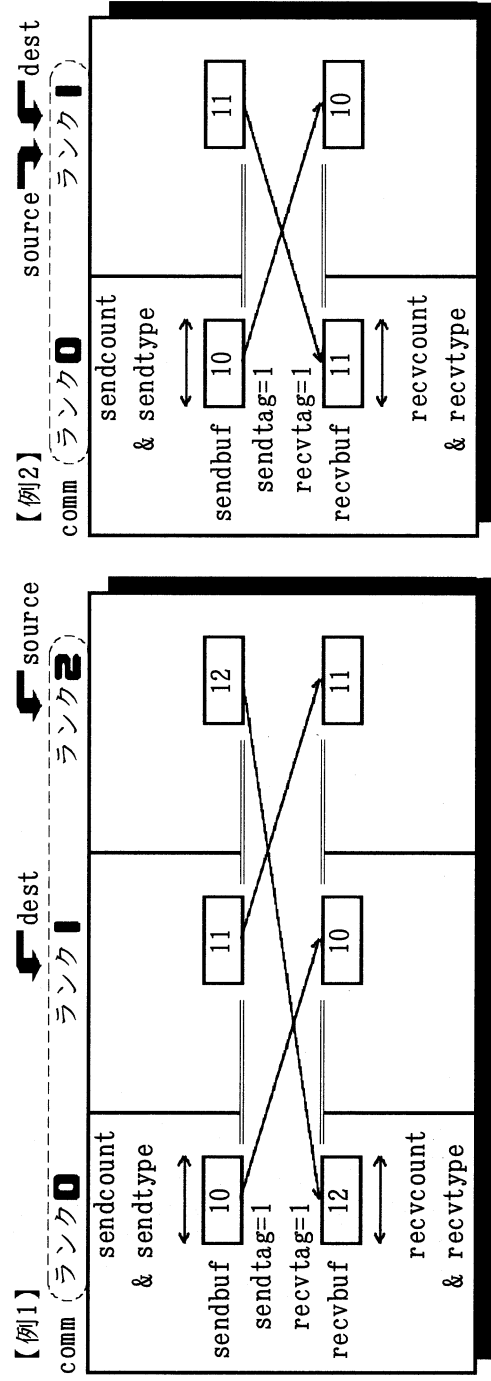
```
[aoyama@node01]~/u/aoyama: mpiexec -np 2 a.out
IRBUF = 9
```

# 1対1ブロッキング通信サブルーチン MPI\_SENDRECV

## 機能 (関連する節: 3-4節, 4-6-2節)

MPI\_SENDRECVは、MPI\_SENDとMPI\_RECVの両機能を1つのサブルーチンで実行します。【例1】のように宛先(dest)と送信元(source)を別にすることも、【例2】のように同一にすることも可能です。

- MPI\_SENDRECVはブロッキング通信サブルーチンですが、デッドロックにはなりません。
- 送信バッファアリアー(sendbuf)内の、データ型がsendtypeであるsendcount個の連続した要素から構成した送信メッセージを、タグ(sendtag)をつけて宛先プロセス(コミュニケータcomm内でランクがdest)に送信します。
- 送信元プロセス(コミュニケータcomm内でランクがsource)から送信された、タグがrecvtagのメッセージを、データ型がrecvtypeであるrecvcount個の連続した要素から構成される受信バッファアリアー(recvbuf)に受信します。
- MPI\_SENDRECVをコールした後、アプリケーションプログラムは、送信バッファアリアーが再利用可能になり、かつ受信バッファアリアーへのメッセージの受信が完了するまで待ちの状態に入ります。



## サンプルプログラム 【例1】のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)
IUP = MOD(MYRANK+1,NPROCS)
IDOWN = MOD(MYRANK-1+NPROCS,NPROCS)
ISBUF = MYRANK + 10
CALL MPI_SENDRECV(ISBUF,1,MPI_INTEGER,IUP,1,IRBUF,1,MPI_INTEGER,IDOWN,1,
&
&
MPI_COMM_WORLD,ISTATUS,IERR)
PRINT *, 'MYRANK = ',MYRANK, ' IRBUF = ',IRBUF
CALL MPI_FINALIZE(IERR)
END
  
```

MPIが使用するインクルードファイルを指定します。  
状況オブジェクトを指定します。  
MPIの環境の初期化処理を行います。  
プロセス数NPROCSを取得します。  
自分のランクMYRANKを取得します。  
宛先のプロセスのランクを指定します。  
送信元のプロセスのランクを指定します。  
各プロセスは送信バッファアリアーにデータをセットします。

```

[aoyama@node01]/u/aoyama: mpi_run_np 3 a.out
MYRANK = 0 IRBUF = 12
MYRANK = 1 IRBUF = 10
MYRANK = 2 IRBUF = 11
  
```

## ■ 使用法

```
CALL MPI_SENDECV(sendbuf, sendcount, sendtype, dest, sendtag,
                 recvbuf, recvcount, recvtype, source, comm, ierror)
```

- **sendbuf** : 送信バッファの先頭アドレスを指定します。
- **sendcount** : 整数。送信バッファ内の送信メッセージの要素数を指定します。
- **sendtype** : 整数。送信バッファ内の送信メッセージのデータ型を指定します。
- **dest** : 整数。宛先プロセスのcomm内でのランクを指定します。宛先プロセスに対してメッセージを送信しない場合には、MPI\_PROC\_NULLを指定します。
- **sendtag** : 整数。送信するメッセージの種類を区別したい場合に使用します。特に区別する必要がない場合には、適当な値(例えば1)を指定します。  
0 ~ MPI\_TAG\_UB までの整数を指定することができます。MPI\_TAG\_UBの値はmpif.hの中で設定されており、マシン環境によって異なります(PRINT文で書き出せば値が分かります)。
- **recvbuf** : 受信バッファの先頭アドレスを指定します。送信バッファと重複してはいけません。
- **recvcount** : 整数。受信バッファの要素数を指定します。受信メッセージの要素数と同じか、それ以上の大きさにして下さい。受信メッセージの要素数より少ない場合は、エラーになります(マシン環境によってはエラーにならない場合もあります)。
- **recvtype** : 整数。受信バッファに入る受信メッセージのデータ型を指定します。
- **source** : 整数。受信したいメッセージを送信するプロセスのcomm内でのランクを指定します。任意のロセスからのメッセージを受信したい場合は、ワイルドカードMPI\_ANY\_SOURCEを指定します。送信元プロセスからのメッセージを受信しない場合には、MPI\_PROC\_NULLを指定します。
- **recvtag** : 整数。受信したいメッセージに付けられているタグの値を指定します。この値は送信元プロセスがそのメッセージを送信した際に引数tagとして付けた値です。
- **comm** : 整数。自分のプロセス、宛先プロセス、送信元プロセスを含むグループのコミュニケーションを指定します。
- **status** : 状況オブジェクトの配列を指定します。サンプルプログラムに示すように、大きさが MPI\_STATUS\_SIZEの整数配列を、『INCLUDE 'mpif.h'』より後に指定して下さい。本サブルーチンが完了すると、受信したメッセージの送信元のランクがstatus(MPI\_SOURCE)に、タグがstatus(MPI\_TAG)に戻ります。これらの値は通常は使用しません。使用するのは、ワイルドカードMPI\_ANY\_SOURCEやMPI\_ANY\_TAGを指定した際、受信したメッセージの送信元のランクやタグの値を知りたい場合です(4-7-1-3節参照)。その場合、受信したメッセージの要素数も知る必要があるならば、『CALL MPI\_GET\_COUNT』を使用して下さい。
- **ierror** : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

## ■ 注意

- **sendcount**, **sendtype**と**recvcount**, **recvtype**は異なっていないともかまいません。

# 1対1非ブロッキング通信サブルーチン MPI\_ISEND

## 機能 (関連する節: 3 - 4 節)

送信バッファ (buf) 内の、データ型が `datatype` である `count` 個の連続した要素から構成した送信メッセージを、タグ (`tag`) をつけて宛先プロセス (コミュニケータ `comm` 内でランクが `dest`) に送信します。MPI\_ISEND で送信したメッセージは、MPI\_RECV と MPI\_Irecv のどちらかで受信してもかまいません。

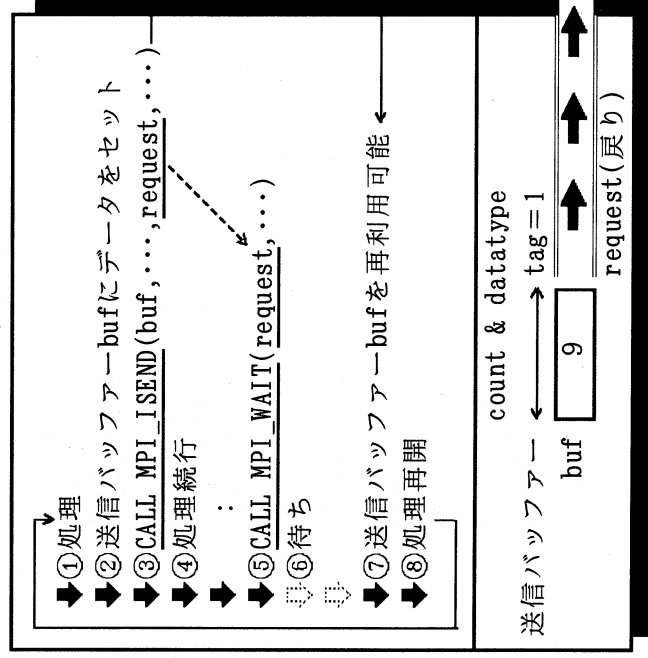
MPI\_ISEND は MPI\_WAIT とペアで使用します。MPI\_ISEND をコールした後、アプリケーションプログラムは、送信バッファが再利用可能になるのを待たず、ただちに次の処理を続行しますので、次に MPI\_WAIT をコールするまでは送信バッファの内容を更新しないで下さい。

MPI\_ISEND を指定して MPI\_WAIT を指定し忘れた場合、タイミングによる (再現性のない) エラーが発生します。また MPI\_ISEND と MPI\_WAIT の引数 `request` がスベルミスで異なっている場合も、当然ながら MPI\_WAIT は働かず、タイミングによる (再現性のない) エラーが発生しますので、スベルミスには十分注意して下さい。

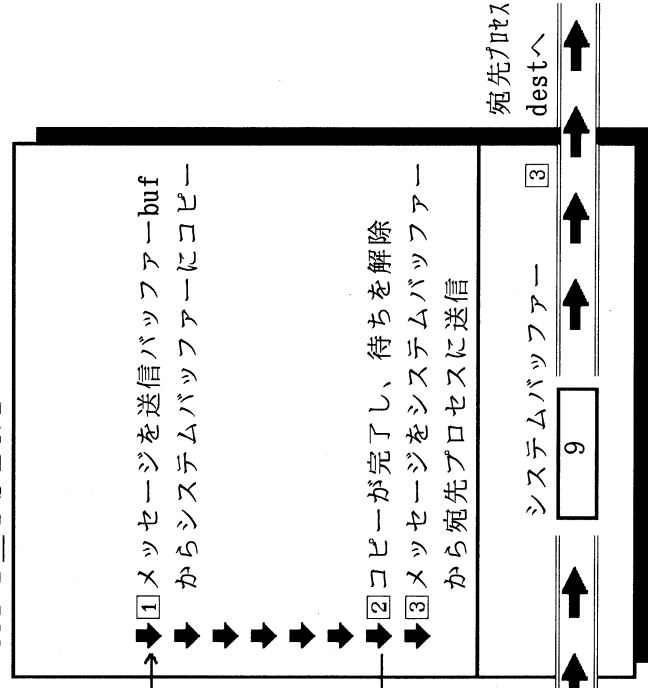
## 動作

- アプリケーションプログラムは、②で送信バッファ `buf` にデータ (メッセージ) をセットし、③で MPI\_ISEND をコールした後、ただちに④で処理を続行します。③の送信メッセージに対し、識別子 (ID) が MPI\_ISEND によって付けられ、引数 `request` に戻ります。
- アプリケーションプログラムは、②で再び送信バッファ `buf` を再利用するよりも前に、⑤で引数 `request` を指定して MPI\_WAIT をコールし、その後⑥で待ちの状態に入ります。
- サブルーチン MPI\_ISEND は、①でメッセージを送信バッファ `buf` からシステムバッファにコピーします。コピーが完了した後、⑤が実行されている場合には②でアプリケーションプログラムの待ちを解除します。なお、待ちが解除されたということは、送信したメッセージが宛先プロセスに到着したことを意味するわけではありません。
- ③の時点で、宛先プロセスが MPI\_RECV または MPI\_Irecv をすでにコールしている場合は、システムバッファのメッセージを宛先プロセスに送信します。宛先プロセスがまだコールしていない場合は、コールするまで③で待ちます。
- アプリケーションプログラムは、⑦で送信バッファ `buf` の再利用が可能になり、⑧で処理を再開します。
- ②で再び送信バッファ `buf` にデータ (メッセージ) をセットしても、前回の②で送信バッファ `buf` にセットしたメッセージはすでに②でシステムバッファにコピーされているので、壊れることはありません。

アプリケーションプログラム



MPI\_ISEND





## ■ 使用法

```
CALL MPI_ISEND(buf, count, datatype, dest, tag, comm, request, ierror)
```

- buf : 送信バッファの先頭アドレスを指定します。
- count : 整数。送信バッファ内の送信メッセージの要素数を指定します。
- datatype : 整数。送信バッファ内の送信メッセージのデータ型を指定します。
- dest : 整数。宛先プロセスのcomm内でのランクを指定します。宛先プロセスに対してメッセージを送信しない場合には、MPI\_PROC\_NULLを指定します。
- tag : 整数。送信するメッセージの種類を区別したい場合に使用します。特に区別する必要がない場合には、適当な値(例えば1)を指定します。
- comm : 0 ~ MPI\_TAG\_UB までの整数を指定することができます。MPI\_TAG\_UBの値はmpif.hの中で設定されており、マシン環境によって異なります(PRINT文で書き出せば値が分かります)。
- request : 整数。自分のプロセスと宛先プロセスを含むグループのコミュニケーションを指定します。
- ierror : 整数。通信リクエスト(送信を要求したメッセージに付けられた識別子)が戻ります。後にMPI\_WAITの引数でこの値を指定します。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

■ サンプルプログラム 【例1】のプログラム例です。ランク0からランク1にメッセージを送信します。

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)  プロセス数NPROCSを取得します。
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)  自分のランクMYRANKを取得します。
  IF (MYRANK == 0) ISBUF = 9
  ランク0のプロセスは送信バッファにデータをセットします。

  IF (MYRANK == 0) THEN
    CALL MPI_ISEND(ISBUF, 1, MPI_INTEGER,  ランク0はランク1にメッセージを送信します。
      1, 1, MPI_COMM_WORLD, IREQ, IERR)
  ELSEIF (MYRANK == 1) THEN
    CALL MPI_Irecv(IRBUF, 1, MPI_INTEGER,  ランク1はランク0からメッセージを受信します。
      0, 1, MPI_COMM_WORLD, IREQ, IERR)
  ENDIF
  CALL MPI_WAIT(IREQ, ISTATUS, IERR)  送信が完了するまで待ちます(受信も同様です)。
  IF (MYRANK == 1) PRINT *, IRBUF  通信後の受信バッファの内容を書き出します。
  CALL MPI_FINALIZE(IERR)
  END

[aoyama@node01]/u/aoyama: mpi_run -np 2 a.out
IRBUF = 9
```

# 1対1非ブロッキング通信サブルーチン MPI\_I\_RECV

## 機能 (関連する節: 3-4 節)

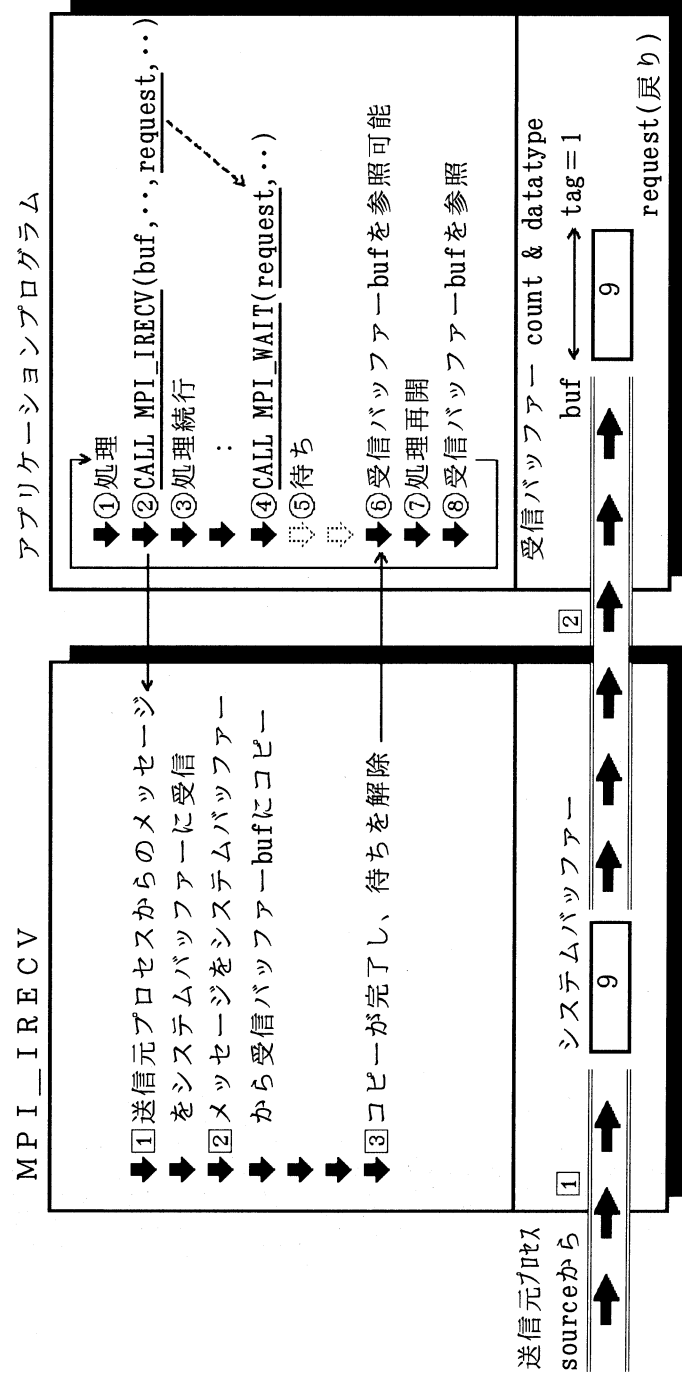
送信元プロセス(コミュニケーションがsource)から送信された、タグがtagのメッセージを、データ型がdatatypeであるcount個の連続した要素から構成される受信バッファア- (buf) に受信します。MPI\_SENDとMPI\_ISENDのどちらで送信したメッセージも、MPI\_I\_RECVで受信することができます。

MPI\_I\_RECVはMPI\_WAITとペアで使用します。MPI\_I\_RECVをコールした後、アプリケーションプログラムは、受信バッファア-へのメッセージの受信が完了するのを待たず、ただちに処理を続行しますので、次にMPI\_WAITをコールするまでは受信バッファア-の内容を参照しないで下さい。

MPI\_RECVを指定してMPI\_WAITを指定し忘れた場合、タイミングによる(再現性のない)エラーが発生します。またMPI\_I\_RECVとMPI\_WAITの引数requestがスペルミスで異なっている場合も、当然ながらMPI\_WAITは働かず、タイミングによる(再現性のない)エラーが発生しますので、スペルミスには十分注意して下さい。

## 動作

- アプリケーションプログラムは、②でMPI\_I\_RECVをコールした後、ただちに③で処理を続行します。②の受信メッセージに対し、識別子(ID)がMPI\_RECVによって付けられ、引数requestに戻ります。
- アプリケーションプログラムは、⑥で受信バッファア-bufを参照するよりも前に、④で引数requestを指定してMPI\_WAITをコールし、その後⑤で待ちの状態に入ります。
- サブルーチンMPI\_RECVは、①の時点で送信元プロセスがすでにMPI\_SENDまたはMPI\_ISENDをコールしている場合は、送信元プロセスから送られたメッセージをシステムバッファア-に受信します。送信元プロセスがまだコールしていない場合は、コールするまで①で待ちます。
- ②でメッセージをシステムバッファア-から受信バッファア-bufにコピーします。コピーが完了した後、④が実行されている場合には③でアプリケーションプログラムの待ちを解除します。
- アプリケーションプログラムは、⑥で受信バッファア-bufを参照可能になり、⑦で処理を再開します。
- ⑧で受信バッファア-bufを参照したときには、受信メッセージが全て入っています。



## 使用法

```
CALL MPI_Irecv(buf, count, datatype, source, tag, comm, request, ierror)
```

- buf : 受信バッファの先頭アドレスを指定します。
- count : 整数。受信バッファの要素数を指定します。受信メッセージの要素数と同じか、それ以上の大きさにして下さい。受信メッセージの要素数より少ない場合は、エラーになります(マシン環境によってはエラーにならない場合もあります)。
- datatype : 整数。受信バッファに入る受信メッセージのデータ型を指定します。
- source : 整数。受信したいメッセージを送信するプロセスのcomm内でのランクを指定します。任意のプロセスからのメッセージを受信したい場合は、ワイルドカードMPI\_ANY\_SOURCEを指定します。送信元プロセスからのメッセージを受信しない場合には、MPI\_PROC\_NULLを指定します。
- tag : 整数。受信したいメッセージに付けられているタグの値を指定します。この値は送信元プロセスがそのメッセージを送信した際に引数tagとして付けた値です。
- comm : 任意のタグ値のメッセージを受信したい場合は、ワイルドカードMPI\_ANY\_TAGを指定します。
- request : 整数。自分のプロセスと送信元プロセスを含むグループのコミュニケーションを指定します。
- ierror : 整数。通信リクエスト(受信を要求したメッセージに付けられた識別子)が戻ります。後にMPI\_WAITの引数でこの値を指定します。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

サンプルプログラム 【例1】 のプログラム例です。ランク 1 はランク 0 からメッセージを受信します。

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER ISTATUS(MPI_STATUS_SIZE)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)  プロセス数NPROCSを取得します。
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)  自分のランクMYRANKを取得します。
  IF (MYRANK == 0) ISBUF = 9
  &
  IF (MYRANK == 0) THEN
    CALL MPI_SEND(ISBUF, 1, MPI_INTEGER, 1, 1, MPI_COMM_WORLD, IREQ, IERR)
    &
    ELSEIF (MYRANK == 1) THEN
      CALL MPI_RECV(IRBUF, 1, MPI_INTEGER, 0, 1, MPI_COMM_WORLD, IREQ, IERR)
    ENDIF
  CALL MPI_WAIT(IREQ, ISTATUS, IERR)  受信が完了するまで待ちます(送信も同様です)。
  IF (MYRANK == 1) PRINT *, IRBUF = ', IRBUF  通信後の受信バッファの内容を書き出します。
  CALL MPI_FINALIZE(IERR)
  END
```

```
[aoyama@node01]/u/aoyama: mpirun -np 2 a.out
IRBUF = 9
```

**機能** (関連する節: 3 - 4 節)

MPI\_WAITは、1対1非ブロッキング通信サブルーチン(MPI\_ISENDまたはMPI\_IRECV)を使用した場合に用いずので、MPI\_ISENDまたはMPI\_IRECVの説明も併せて参照して下さい。

MPI\_ISENDまたはMPI\_IRECVを指定してMPI\_WAITを指定し忘れた場合、タイミングによる(再現性のない)エラーが発生します。またMPI\_ISENDまたはMPI\_IRECVとMPI\_WAITの引数requestがスペルミスで異なっている場合も、当然ながらMPI\_WAITは働かず、タイミングによる(再現性のない)エラーが発生しますので、スペルミスには十分注意して下さい。

**使用法**

```
CALL MPI_WAIT(request, status, ierror)
```

- request : 整数。通信リクエストを指定します。これはMPI\_ISENDまたはMPI\_IRECVをコールしたときにメッセージに付けられた識別子requestです。  
この引数の内容は、MPI\_WAITが完了するとMPI\_REQUEST\_NULLという値に変わります。
- status : 状況オブジェクトの配列を指定します。サンプルプログラムに示すように、大きさがMPI\_STATUS\_SIZEの整数配列を、『INCLUDE 'mpif.h'』より後に指定して下さい。  
MPI\_IRECVに対してMPI\_WAITをコールした場合、MPI\_WAITが完了すると、受信したメッセージの送信元のランクの値がstatus(MPI\_SOURCE)に、タグの値がstatus(MPI\_TAG)に戻ります(サンプルプログラム参照)。  
これらの値は通常のプログラムでは使用しません。使用するのは、サンプルプログラムのように、MPI\_IRECVの引数にファイルドカードMPI\_ANY\_SOURCEやMPI\_ANY\_TAGを指定し、受信した後にメッセージの送信元のランクやタグの値を知りたい場合です(4-7-1-3節参照)。  
その場合、受信したメッセージの要素数も知る必要があるならば、『CALL MPI\_GET\_COUNT』を使用して下さい。
- MPI\_ISENDに対してMPI\_WAITをコールした場合、MPI\_WAITが完了するとstatus(MPI\_SOURCE)に値がセットされますが、この値は処理系によって異なり、またプログラムで使用することはありません。ただし、値がセットされるので配列status(MPI\_SOURCE)の宣言は必要です。  
また、4-7-1-3節のマスターズレスレブ方式などでMPI\_ISENDとMPI\_IRECVの両方を使用して、前述のようにMPI\_IRECVに対するMPI\_WAITでセットされるstatusの値をプログラムで使用したい場合、MPI\_ISENDに対するMPI\_WAITでセットされるstatusと別の配列名にしないと値が壊れてしまいますので注意して下さい。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

サンプルプログラム ランク0からランク1にメッセージを送信します。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。
IF (MYRANK == 0) THEN
  ISTATUS = 9
ENDIF
IF (MYRANK == 0) THEN
  CALL MPI_ISEND(ISBUF, 1, MPI_INTEGER, ランク0はランク1にメッセージを送信します。
    1, 1, MPI_COMM_WORLD, IREQ, IERR)
ELSEIF (MYRANK == 1) THEN
  CALL MPI_IRECV(IRBUF, 1, MPI_INTEGER, ランク1は任意のプロセスからの任意のタグの
    MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, IREQ, IERR)
ENDIF
CALL MPI_WAIT(IREQ, ISTATUS, IERR) 送信(受信)が完了するまで待ちます。
IF (MYRANK == 1) THEN
  PRINT *, 'IRBUF = ', IRBUF 通信後の受信バッファの内容を書き出します。
  PRINT *, 'SOURCE = ', ISTATUS(MPI_SOURCE) 送信元プロセスのランクを書き出します。
  PRINT *, 'TAG = ', ISTATUS(MPI_TAG) 受信したメッセージのタグを書き出します。
ENDIF
CALL MPI_FINALIZE(IERR)
END

```

```

[aoyama@node01]/u/aoyama: mpiexec -np 2 a.out
IRBUF = 9
SOURCE = 0
TAG = 1

```

## 1対1通信サブルーチン MPI\_GET\_COUNT

### 機能 (関連する節: 4-6-10 節)

1対1通信サブルーチンMPI\_RECV(またはMPI\_IRecv)によって受信したメッセージの要素数を調べたい場合に本サブルーチンを使用します。MPI\_RECV(またはMPI\_IRecv)に対応するMPI\_WAITが完了した後、MPI\_RECV(またはMPI\_WAIT)の引数で指定した状況オブジェクトstatusに各種の情報が戻ります。このstatusと、MPI\_RECV(またはMPI\_IRecv)の引数で指定した受信メッセージのデータ型(datatype)を引数で指定して本サブルーチンをコールすると、MPI\_RECV(またはMPI\_IRecv)で受信したメッセージの要素数がcountに戻ります。

### 使用法

```
CALL MPI_GET_COUNT(status, datatype, count, ierror)
```

- status : 大きさがMPI\_STATUS\_SIZEの整数配列。調べたいMPI\_RECV(またはMPI\_IRecv)に対応するMPI\_WAITの引数で指定した状況オブジェクトstatusを指定します。
- datatype : 整数。MPI\_RECV(またはMPI\_IRecv)の引数で指定した受信メッセージのデータ型(datatype)を指定します。
- count : 整数。MPI\_RECV(またはMPI\_IRecv)によって受信したメッセージの要素数が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

### 注意

- 本サブルーチンは、調べたいMPI\_RECV(またはMPI\_IRecv)に対応するMPI\_WAITが完了した後でコールする必要があります。

サンプルプログラム ランク0からランク1にメッセージを送信します。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER ISTATUS(MPI_STATUS_SIZE)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR)  プロセス数NPROCSを取得します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR)  自分のランクMYRANKを取得します。
IF (MYRANK == 0) THEN
  ISBUF = 9
  ENDF
IF (MYRANK == 0) THEN
  CALL MPI_ISEND(ISBUF,1,MPI_INTEGER, ランク0はランク1にメッセージを送信します。
    1,1,MPI_COMM_WORLD,IREQ,IERR)
ELSEIF (MYRANK == 1) THEN
  CALL MPI_IRECV(IRBUF,1,MPI_INTEGER, ランク1は任意のプロセスからの任意のタグの
    MPI_ANY_SOURCE,MPI_ANY_TAG,MPI_COMM_WORLD,IREQ,IERR)  メッセージを受信します。
  ENDF
CALL MPI_WAIT(IREQ,ISTATUS,IERR)  送信(受信)が完了するまで待ちます。
CALL MPI_GET_COUNT(ISTATUS,MPI_INTEGER,ICOUNT,IERR)  受信したメッセージの要素数を調べます。
IF (MYRANK == 1) THEN
  PRINT *, 'IRBUF = ',IRBUF  通信後の受信バッファの内容を書き出します。
  PRINT *, 'SOURCE = ',ISTATUS(MPI_SOURCE)  送信元プロセスのランクを書き出します。
  PRINT *, 'TAG = ',ISTATUS(MPI_TAG)  受信したメッセージのタグを書き出します。
  PRINT *, 'COUNT = ',ICOUNT  受信したメッセージの要素数を書き出します。
  ENDF
CALL MPI_FINALIZE(IERR)
END

```

```
[aoyama@node01]/u/aoyama: mpirun -np 2 a.out
```

```

IRBUF = 9
SOURCE = 0
TAG = 1
COUNT = 1

```

# 派生データ型に関するサブルーチン MPI\_TYPE\_SIZE

## 機能 (関連する節: 3 - 5 - 2 節)

指定したデータ型 (datatype) の大きさが size に戻ります。C 言語の sizeof 関数に相当します。

## 使用法

```
CALL MPI_TYPE_SIZE(datatype, size, ierror)
```

- datatype : 整数。大きさを調べたいデータ型を指定します。
- size : 整数。指定したデータ型 (datatype) の大きさ (単位は バイト) が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

### 【例1】

```
datatype = MPI_INTEGER
```

size には 4 (バイト) が戻ります。

### 【例2】

```
datatype = MPI_INTEGER MPI_INTEGER MPI_INTEGER
```

size には 12 (バイト) が戻ります。

## サンプルプログラム 【例1】 のプログラム例です。

```
PROGRAM MAIN
  INCLUDE 'mpif.h'
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
  CALL MPI_TYPE_SIZE(MPI_INTEGER, ISIZE, IERR)
  IF (MYRANK=0) PRINT *, 'ISIZE = ', ISIZE
  CALL MPI_FINALIZE(IERR)
END
```

MPI が使用するインクルードファイルを指定します。  
MPI 環境の初期化処理を行います。  
プロセス数 NPROCS を取得します。  
自分のランク MYRANK を取得します。  
MPI\_INTEGER の大きさを調べます。  
結果を書き出します。  
MPI 環境の終了処理を行います。

```
[aoyama@node01]/u/aoyama: mpiexec -np 1 a.out
```

```
ISIZE = 4
```



機能 (関連する節: 3-5-2 節)

- count個の古いデータ型(oldtype)を連結して新しいデータ型(newtype)を作成します。
- 作成した新しいデータ型(newtype)は、MPI\_TYPE\_COMMITで登録してから通信ルーチンで使用して下さい。

使用法

```
CALL MPI_TYPE_CONTIGUOUS(count, oldtype, newtype, ierror)
```

- count : 整数。古いデータ型(oldtype)をいくつ連結するかを指定します。
- oldtype : 整数。古いデータ型を指定します。
- newtype : 整数。作成された新しいデータ型が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

【例1】

```
1 2 3 = count
newtype = [●][●][●] ● : MPI_INTEGER(oldtype)
```

サンプルプログラム 【例1】のプログラム例です。

```
PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER IBUF(10)

CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。

IF (MYRANK == 0) THEN
  DO I = 1, 10
    IBUF(I) = I
  ENDDO
ENDIF

CALL MPI_TYPE_CONTIGUOUS(3, MPI_INTEGER, INEWTYPE, IERR) 新しいデータ型を作成します。
CALL MPI_TYPE_COMMIT(INEWTYPE, IERR) 新しいデータ型を登録します。

CALL MPI_BCAST(IBUF, 1, INEWTYPE, 0, MPI_COMM_WORLD, IERR) 新しいデータ型でメッセージを
PRINT *, 'MYRANK = ', MYRANK, 'IBUF = ', IBUF 1個通信し、書き出します。
CALL MPI_FINALIZE(IERR) MPI環境の終了処理を行います。
END
```

```
[aoyama@node01]/u/aoyama: mpirun -np 2 a.out
MYRANK = 0 IBUF = 1 2 3 4 5 6 7 8 9 10
MYRANK = 1 IBUF = 1 2 3 0 0 0 0 0 0
```

## 機能 (関連する節: 3-5-2 節)

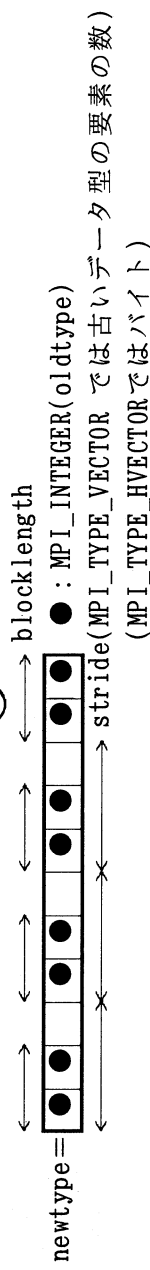
- count個のブロックから構成される新しいデータ型(newtype)を作成します。
- 各ブロックには、古いデータ型(oldtype)の要素がblocklength個含まれます。
- 各ブロックの先頭間の間隔はstrideです。
- 作成した新しいデータ型(newtype)は、MPI\_TYPE\_COMMITで登録してから通信ルーチンで使用して下さい。

## 使用法

```
CALL MPI_TYPE_VECTOR (count, blocklength, stride, oldtype, newtype, ierror)
CALL MPI_TYPE_HVECTOR(count, blocklength, stride, oldtype, newtype, ierror)
```

- count : 整数。ブロックの数を指定します。
- blocklength : 整数。各ブロックに含まれる古いデータ型の要素の数(全ブロックで同一)を指定します。
- stride : 整数。各ブロックの先頭間の間隔(一定)を指定します。  
単位は、MPI\_TYPE\_VECTOR では古いデータ型の要素の数、MPI\_TYPE\_HVECTORではバイトです。
- oldtype : 整数。古いデータ型(全ブロックで同一)を指定します。
- newtype : 整数。作成された新しいデータ型が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

【例1】 ブロック 1 2 3 ④=count



## サンプルプログラム 【例1】のプログラム例です。

```

PROGRAM MAIN
  INCLUDE 'mpif.h'
  INTEGER IBUF1(20), IBUF2(20)
  CALL MPI_INIT(IERR)
  CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)  プロセス数NPROCSを取得します。
  CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)  自分のランクMYRANKを取得します。
  IF (MYRANK == 0) THEN
    DO I = 1, 20
      IBUF1(I) = I
      IBUF2(I) = I
    ENDDO
  ELSE
    DO I = 1, 20
      IBUF1(I) = 0
      IBUF2(I) = 0
    ENDDO
  ENDIF
  CALL MPI_TYPE_VECTOR(4, 2, 3, MPI_INTEGER, INEWTYPE1, IERR)  新しいデータ型を作成します。
  CALL MPI_TYPE_COMMIT(INEWTYPE1, IERR)  新しいデータ型を登録します。
  CALL MPI_BCAST(IBUF1, 1, INEWTYPE1, 0, MPI_COMM_WORLD, IERR)  新しいデータ型でメッセージを
  PRINT *, 'MYRANK = ', MYRANK, 'IBUF1 = ', IBUF1  1個通信し、書き出します。
  CALL MPI_TYPE_HVECTOR(4, 2, 3*4, MPI_INTEGER, INEWTYPE2, IERR)  新しいデータ型を作成します。
  CALL MPI_TYPE_COMMIT(INEWTYPE2, IERR)  新しいデータ型を登録します。
  CALL MPI_BCAST(IBUF2, 1, INEWTYPE2, 0, MPI_COMM_WORLD, IERR)  新しいデータ型でメッセージを
  PRINT *, 'MYRANK = ', MYRANK, 'IBUF2 = ', IBUF2  1個通信し、書き出します。
  CALL MPI_FINALIZE(IERR)
  END

```

```

[aoyama@node01]/u/aoyama: mpirun -np 2 a.out
MYRANK = 0 IBUF1 = 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
MYRANK = 1 IBUF1 = 1 2 0 4 5 0 7 8 0 10 11 0 0 0 0 0 0 0 0
MYRANK = 0 IBUF2 = 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
MYRANK = 1 IBUF2 = 1 2 0 4 5 0 7 8 0 10 11 0 0 0 0 0 0 0 0

```

## 【MPI-2】

### 派生データ型に関するサブルーチン MPI\_TYPE\_CREATE\_INDEXED\_BLOCK

■ 機能 (関連する節: 3-5-2 節, 4-6-6-1 節)

- 本ルーチンはMPI-2のサブルーチンです。
- count個のブロックから構成される新しいデータ型(newtype)を作成します。
- 各ブロックには、古いデータ型(oldtype)の要素がblocklength個含まれます。
- ブロックごとに、新しいデータ型の先頭から各ブロックの先頭までの変位(array\_of\_displacements)を指定します。
- 作成した新しいデータ型(newtype)は、MPI\_TYPE\_COMMITで登録してから通信ルーチンで使用して下さい。

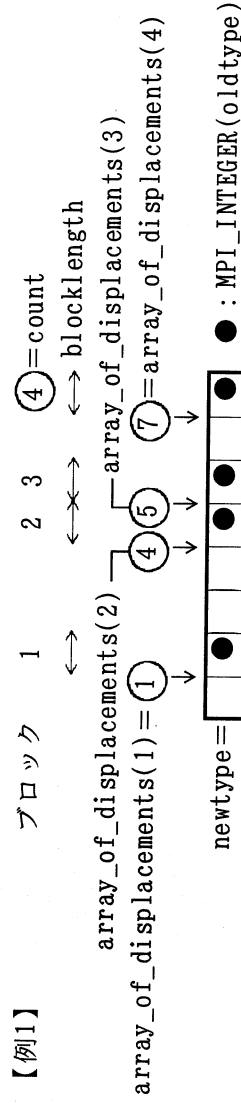
#### ■ 使用法

```
CALL MPI_TYPE_CREATE_INDEXED_BLOCK(count, blocklength, array_of_displacements,  
oldtype, newtype, ierror)
```

- count : 整数。ブロックの数を指定します。  
この値は同時に配列array\_of\_displacementsの要素数を表します。
- blocklength : 整数。各ブロックに含まれる古いデータ型の要素の数(全ブロックで同一)を指定します。
- array\_of\_displacements : 整数配列。新しいデータ型の先頭からi個目のブロックの先頭までの変位をarray\_of\_displacements(i)に指定します。単位は古いデータ型の要素の数です。
- oldtype : 整数。古いデータ型(全ブロックで同一)を指定します。
- newtype : 整数。作成された新しいデータ型が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

#### ■ 注意

最初のブロックの前にブランクがあっても構いません。ただし、この派生データ型を使用してデータを2個以上通信した場合や、この派生データ型をMPI\_TYPE\_CONTIGUOUSで2個以上連続させた場合、2個目以降の派生データ型では前のブランクが取れてしまいます。ブランクを残したい場合は、3-5-2 節の(7)~(9)を参照して下さい。



## ■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER IBUF(10)
INTEGER IDISP(4)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
IF (MYRANK == 0) THEN
  DO I = 1, 10
    IBUF(I) = I
  ENDDO
ELSE
  DO I = 1, 10
    IBUF(I) = 0
  ENDDO
ENDIF
IDISP(1) = 1
IDISP(2) = 4
IDISP(3) = 5
IDISP(4) = 7
CALL MPI_TYPE_CREATE_INDEXED_BLOCK
& (4, 1, IDISP, MPI_INTEGER, INEWTYPE, IERR)
CALL MPI_TYPE_COMMIT(INEWTYPE, IERR)
CALL MPI_BCAST(IBUF, 1, INEWTYPE, 0, MPI_COMM_WORLD, IERR)
PRINT *, 'MYRANK = ', MYRANK, 'IBUF = ', IBUF
CALL MPI_FINALIZE(IERR)
END

```

MPIが使用するインクルードファイルを指定します。  
送(受)信バッファを指定します。

MPI\_TYPE\_CREATE\_INDEXED\_BLOCKで使用する引数の配列を指定します。

MPI環境の初期化処理を行います。

CALL MPI\_COMM\_SIZE(MPI\_COMM\_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。

CALL MPI\_COMM\_RANK(MPI\_COMM\_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。

ランク0のプロセスは送信バッファに  
データをセットします

ランク0以外のプロセスは受信バッファを  
ゼロクリアします。

新しいデータ型を作成するために

MPI\_TYPE\_CREATE\_INDEXED\_BLOCKの各引数を指定します。

新しいデータ型を作成します。

CALL MPI\_TYPE\_CREATE\_INDEXED\_BLOCK

& (4, 1, IDISP, MPI\_INTEGER, INEWTYPE, IERR)

CALL MPI\_TYPE\_COMMIT(INEWTYPE, IERR)

CALL MPI\_BCAST(IBUF, 1, INEWTYPE, 0, MPI\_COMM\_WORLD, IERR)

PRINT \*, 'MYRANK = ', MYRANK, 'IBUF = ', IBUF

CALL MPI\_FINALIZE(IERR)

END

[aoyama@node01]/u/aoyama: mpiexec -np 2 a.out

MYRANK = 0 IBUF = 1 2 3 4 5 6 7 8 9 10

MYRANK = 1 IBUF = 0 2 0 0 5 6 0 8 0 0



## サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER IBUF1(10),IBUF2(10)
INTEGER IBLOCK(2),IDISP(2)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR) プロセス数NPROCSを取得します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。
IF (MYRANK == 0) THEN
  DO I = 1, 10
    IBUF1(I) = I
    IBUF2(I) = I
  ENDDO
ELSE
  DO I = 1, 10
    IBUF1(I) = 0
    IBUF2(I) = 0
  ENDDO
ENDIF
IBLOCK(1) = 2
IBLOCK(2) = 3
IDISP(1) = 0
IDISP(2) = 3
CALL MPI_TYPE_INDEXED(2,IBLOCK,IDISP,MPI_INTEGER,INETYPE1,IERR) 新しいデータ型を作成します。
CALL MPI_TYPE_COMMIT(INETYPE1,IERR) 新しいデータ型を登録します。
CALL MPI_BCAST(IBUF1,1,INETYPE1,0,MPI_COMM_WORLD,IERR) 新しいデータ型でメッセージ
PRINT *, 'MYRANK = ',MYRANK, 'IBUF1 = ',IBUF1
IBLOCK(1) = 2
IBLOCK(2) = 3
IDISP(1) = 0*4
IDISP(2) = 3*4
CALL MPI_TYPE_HINDEXED(2,IBLOCK,IDISP,MPI_INTEGER,INETYPE2,IERR) 新しいデータ型を作成します。
CALL MPI_TYPE_COMMIT(INETYPE2,IERR) 新しいデータ型を登録します。
CALL MPI_BCAST(IBUF2,1,INETYPE2,0,MPI_COMM_WORLD,IERR) 新しいデータ型でメッセージ
PRINT *, 'MYRANK = ',MYRANK, 'IBUF2 = ',IBUF2
CALL MPI_FINALIZE(IERR)
END

```

```

[aoyama@node01]~/u/aoyama: mpirun -np 2 a.out [J]
MYRANK = 0 IBUF1 = 1 2 3 4 5 6 7 8 9 10
MYRANK = 1 IBUF1 = 1 2 0 4 5 6 0 0 0 0
MYRANK = 0 IBUF2 = 1 2 3 4 5 6 7 8 9 10
MYRANK = 1 IBUF2 = 1 2 0 4 5 6 0 0 0 0

```





## サンプルプログラム 【例1】と【例2】のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
INTEGER IBUF1(10),IBUF2(10)
INTEGER IBLOCK(2),IDISP(2),ITYPE(2)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD,NPROCS,IERR) プロセス数NPROCSを取得します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。
IF (MYRANK == 0) THEN
  DO I = 1, 10
    IBUF1(I) = I
    IBUF2(I) = I
  ENDDO
ELSE
  DO I = 1, 10
    IBUF1(I) = 0
    IBUF2(I) = 0
  ENDDO
ENDIF
IBLOCK(1) = 2
IBLOCK(2) = 3
IDISP(1) = 0*4
IDISP(2) = 3*4
ITYPE(1) = MPI_INTEGER
ITYPE(2) = MPI_INTEGER

CALL MPI_TYPE_STRUCT(2,IBLOCK,IDISP,ITYPE,INETYPE1,IERR) 【例1】の新データ型を作成します。
CALL MPI_TYPE_COMMIT(INETYPE1,IERR) 【例1】の新データ型を登録します。
CALL MPI_BCAST(IBUF1,1,INETYPE1,0,MPI_COMM_WORLD,IERR) 【例1】の新データ型でメッセージを
PRINT *,MYRANK = ',MYRANK, 'IBUF1 = ',IBUF1
IBLOCK(1) = 1
IBLOCK(2) = 3
IDISP(1) = 0*4
IDISP(2) = 2*4
ITYPE(1) = MPI_LB
ITYPE(2) = MPI_INTEGER

CALL MPI_TYPE_STRUCT(2,IBLOCK,IDISP,ITYPE,INETYPE2,IERR) 【例2】の新データ型を作成します。
CALL MPI_TYPE_COMMIT(INETYPE2,IERR) 【例2】の新データ型を登録します。
CALL MPI_BCAST(IBUF2,2,INETYPE2,0,MPI_COMM_WORLD,IERR) 【例2】の新データ型でメッセージを
PRINT *,MYRANK = ',MYRANK, 'IBUF2 = ',IBUF2
CALL MPI_FINALIZE(IERR)
END

```

MPIが使用するインクルードファイルを指定します。  
送(受)信バッファを指定します。

MPI\_TYPE\_STRUCTで使用する引数の配列を指定します。  
MPI環境の初期化処理を行います。

CALL MPI\_COMM\_SIZE(MPI\_COMM\_WORLD,NPROCS,IERR) プロセス数NPROCSを取得します。  
CALL MPI\_COMM\_RANK(MPI\_COMM\_WORLD,MYRANK,IERR) 自分のランクMYRANKを取得します。

IF (MYRANK == 0) THEN

DO I = 1, 10  
IBUF1(I) = I  
IBUF2(I) = I  
ENDDO

【0】のプロセスは送信バッファに  
データをセットします

ELSE  
DO I = 1, 10  
IBUF1(I) = 0  
IBUF2(I) = 0  
ENDDO

【0】以外のプロセスは受信バッファを  
ゼロクリアします。

ENDIF

IBLOCK(1) = 2  
IBLOCK(2) = 3

IDISP(1) = 0\*4  
IDISP(2) = 3\*4

【例1】の新しいデータ型を作成するために  
MPI\_TYPE\_STRUCTの各引数を指定します。

ITYPE(1) = MPI\_INTEGER  
ITYPE(2) = MPI\_INTEGER

CALL MPI\_TYPE\_STRUCT(2,IBLOCK,IDISP,ITYPE,INETYPE1,IERR) 【例1】の新データ型を作成します。  
CALL MPI\_TYPE\_COMMIT(INETYPE1,IERR) 【例1】の新データ型を登録します。

CALL MPI\_BCAST(IBUF1,1,INETYPE1,0,MPI\_COMM\_WORLD,IERR) 【例1】の新データ型でメッセージを  
PRINT \*,MYRANK = ',MYRANK, 'IBUF1 = ',IBUF1  
1個通信し、書き出します。

IBLOCK(1) = 1  
IBLOCK(2) = 3

IDISP(1) = 0\*4  
IDISP(2) = 2\*4

【例2】の新しいデータ型を作成するために  
MPI\_TYPE\_STRUCTの各引数を指定します。

ITYPE(1) = MPI\_LB  
ITYPE(2) = MPI\_INTEGER

CALL MPI\_TYPE\_STRUCT(2,IBLOCK,IDISP,ITYPE,INETYPE2,IERR) 【例2】の新データ型を作成します。  
CALL MPI\_TYPE\_COMMIT(INETYPE2,IERR) 【例2】の新データ型を登録します。

CALL MPI\_BCAST(IBUF2,2,INETYPE2,0,MPI\_COMM\_WORLD,IERR) 【例2】の新データ型でメッセージを  
PRINT \*,MYRANK = ',MYRANK, 'IBUF2 = ',IBUF2  
2個通信し、書き出します。

CALL MPI\_FINALIZE(IERR)  
MPI環境の終了処理を行います。

END

```

[aoyama@node01]/u/aoyama: mpirun -np 2 a.out
MYRANK = 0 IBUF1 = 1 2 3 4 5 6 7 8 9 10 【例1】の結果です。
MYRANK = 1 IBUF1 = 1 2 0 4 5 6 0 0 0 0
MYRANK = 0 IBUF2 = 1 2 3 4 5 6 7 8 9 10 【例2】の結果です。
MYRANK = 1 IBUF2 = 0 0 3 4 5 0 0 8 9 10

```

# 【MPI-2】派生データ型に関するサブルーチン MPI\_TYPE\_CREATE\_SUBARRAY

機能 (関連する節: 3-5-3-1節、4-6-3-1節、4-6-3-2節)

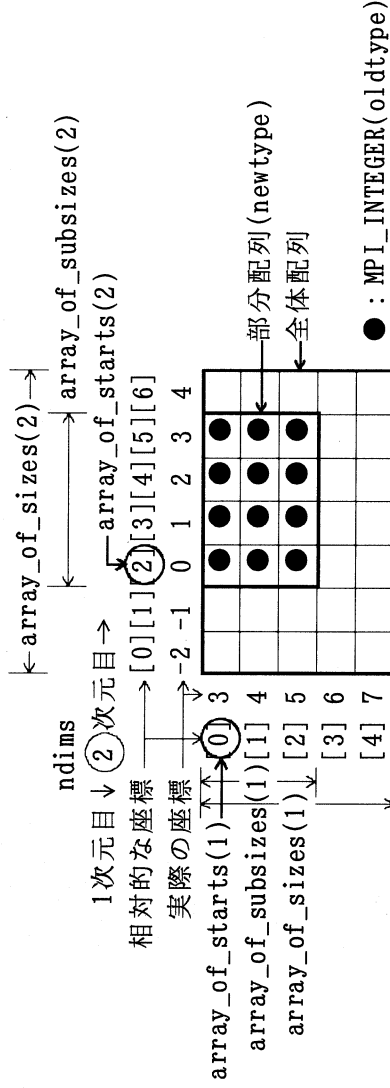
- 【例1】に示すように、全体配列内に含まれる部分配列を表す新しいデータ型(newtype)を作成します。配列の次元数(ndims)は何次元でも可能です。
- 作成した新しいデータ型(newtype)は、MPI\_TYPE\_COMMITで登録してから通信ルーチンで使用して下さい。
- 作成した新しいデータ型(newtype)を通信ルーチンで使用する際、送/受信バッファの先頭アドレスの引数には、全体配列の先頭アドレス(全体配列名でも可)を指定します。
- 3-5-3節で紹介した派生データ型のユーザーリテイナー・サブルーチン PARA\_TYPE\_BLOCK2、PARA\_TYPE\_BLOCK3は、本ルーチンと同様の機能を持ちます。本ルーチンはMPI-2に含まれているので、MPI-2が使用できない環境の場合、代替として3-5-3節のルーチンを使用して下さい。

## 使用法

```
CALL MPI_TYPE_CREATE_SUBARRAY(ndims, array_of_sizes, array_of_subsizes, array_of_starts,  
order, oldtype, newtype, ierror)
```

- ndims : 整数。全体配列の次元数を指定します。
- array\_of\_sizes : 大きさndimsの整数配列。各次元ごとに、全体配列の古いデータ型の要素の数を指定します。
- array\_of\_subsizes : 大きさndimsの整数配列。各次元ごとに、部分配列の古いデータ型の要素の数を指定します。
- array\_of\_starts : 大きさndimsの整数配列。各次元ごとに、部分配列の最初の要素の座標を指定します。座標は【例1】の[\_]に示す相対的な座標(全体配列の最初の要素の座標を0としたときの座標)で指定します。
- order : 整数。Fortranの場合はMPI\_ORDER\_FORTRANを、Cの場合はMPI\_ORDER\_Cを指定します。
- oldtype : 整数。古いデータ型を指定します。
- newtype : 整数。作成された新しいデータ型が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

### 【例1】



## ■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
INCLUDE(mpi.f.h)
INTEGER IBUF(3:7,-2:4)
INTEGER ISIZE(2), ISUBSIZE(2), ISTART(2)
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR) プロセス数NPROCSを取得します。
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR) 自分のランクMYRANKを取得します。
IF(MYRANK==0) THEN
  DO J=-2,4
  DO I=3,7
    IBUF(I,J) = I*10+J
  ENDDO
  ENDDO
ELSE
  DO J=-2,4
  DO I=3,7
    IBUF(I,J) = 0
  ENDDO
  ENDDO
ENDIF
ISIZE(1)=5; ISUBSIZE(1)=3; ISTART(1)=0 新しいデータ型(部分配列)を作成するために、
ISIZE(2)=7; ISUBSIZE(2)=4; ISTART(2)=2 MPI_TYPE_CREATE_SUBARRAYの各引数を指定します。
CALL MPI_TYPE_CREATE_SUBARRAY(2, ISIZE, ISUBSIZE, ISTART, 新しいデータ型を作成します。
& MPI_ORDER_FORTRAN, MPI_INTEGER, INEWTYPE, IERR)
CALL MPI_TYPE_COMMIT(INEWTYPE, IERR) 新しいデータ型を登録します。
CALL MPI_BCAST(IBUF, 1, INEWTYPE, 0, MPI_COMM_WORLD, IERR) 新しいデータ型でメッセージを
IF (MYRANK==1) PRINT *, IBUF = ', IBUF
CALL MPI_FINALIZE(IERR) MPI環境の終了処理を行います。
END

```

```

[aoyama@node01]/u/aoyama: mpirun -np 2 a.out
IBUF = 0 0 0 0 0 0 0 0 30 40 50 0 0 31 41 51 0 0 32 42 52 0 0 33 43 53 0 0 0 0 0 0

```

## 派生データ型に関するサブルーチン MPI\_TYPE\_COMMIT

■ 機能 (関連する節: 3-5-2 節)

MPI\_TYPE\_CONTIGUOUS、MPI\_TYPE\_CREATE\_INDEXED\_BLOCK、MPI\_TYPE\_(H)INDEXED、MPI\_TYPE\_STRUCT、MPI\_TYPE\_CREATE\_SUBARRAYなどで作成した派生データ型(datatype)は、通信ルーチンで使用する前に、MPI\_TYPE\_COMMITで一度だけ登録します。登録が完了すると、そのデータ型は通信ルーチンのデータ型引数に何度でも指定することができます。

■ 使用法

```
CALL MPI_TYPE_COMMIT(datatype, ierror)
```

- datatype : 整数。登録したい、MPI\_TYPE\_CONTIGUOUS、MPI\_TYPE\_(H)VECTOR、MPI\_TYPE\_CREATE\_INDEXED\_BLOCK、MPI\_TYPE\_(H)INDEXED、MPI\_TYPE\_STRUCT、MPI\_TYPE\_CREATE\_SUBARRAYなどで作成した派生データ型を指定します。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

■ 注意

- 新しいデータ型が、通信ルーチンで実際に使用する派生データ型でなく、途中で作られる一時的な派生データ型の場合、本サブルーチンで登録する必要はありません。

■ サンプルプログラム

MPI\_TYPE\_CONTIGUOUS、MPI\_TYPE\_(H)VECTOR、MPI\_TYPE\_CREATE\_INDEXED\_BLOCK、MPI\_TYPE\_(H)INDEXED、MPI\_TYPE\_STRUCT、MPI\_TYPE\_CREATE\_SUBARRAYなどの例を参照して下さい。

### 機能 (関連する節: 3-5-2 節)

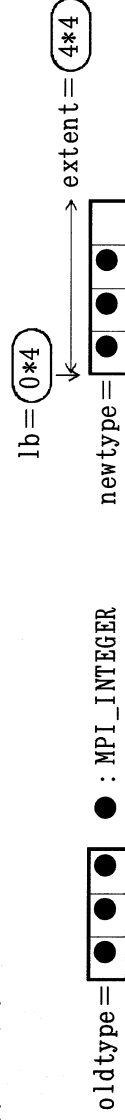
- 本ルーチンはMPI-2のサブルーチンです。
- 本サブルーチンは、古いデータ型(oldtype)の後にブランクを入れるときに使用します(3-5-2 節の(7)~(9)参照)。
- 古いデータ型(oldtype)に対して下限の位置(lb)と範囲(extent)を設定し、新しい(派生)データ型(newtype)を作成します。上限の位置は「lb + extent」になります。
- 作成した新しいデータ型(newtype)は、MPI\_TYPE\_COMMITで登録してから通信ルーチンで使用して下さい。

### 使用法

```
CALL MPI_TYPE_CREATE_RESIZED(oldtype, lb, extent, newtype, ierror)
```

- oldtype : 整数。古いデータ型を指定します。
- lb : 整数(INTEGER(KIND=MPI\_ADDRESS\_KIND)で宣言して下さい)。新しいデータ型の下限の位置(先頭からの変位)を指定します(単位はバイト)。
- extent : 整数(INTEGER(KIND=MPI\_ADDRESS\_KIND)で宣言して下さい)。新しいデータ型の範囲を指定します(単位はバイト)。
- newtype : 整数。作成された新しいデータ型が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

【例1】 付録のMPI\_TYPE\_CONTIGUOUSの【例1】で作成した派生データ型(以下ではIOLDTYPE)の後ろにブランクを1つ付けます。



サンプルプログラム 【例1】のプログラム例です(付録のMPI\_TYPE\_CONTIGUOUSの【例1】参照)。

```

:
INTEGER(KIND=MPI_ADDRESS_KIND) ILB, IEXTENT
:
CALL MPI_TYPE_CONTIGUOUS(3, MPI_INTEGER, IOLDTYPE, IERR)
ILB = 0*4
IEXTENT = 4*4
CALL MPI_TYPE_CREATE_RESIZED(IOLDTYPE, ILB, IEXTENT, INEWTYPE, IERR)
CALL MPI_TYPE_COMMIT(INEWTYPE, IERR)
CALL MPI_BCAST(IBUF, 2, INEWTYPE, 0, MPI_COMM_WORLD, IERR)
PRINT *, 'MYRANK = ', MYRANK, 'IBUF = ', IBUF
:

```

```

[ aoyama@node01 ] /u/aoyama: mpiexec -np 2 a.out
MYRANK = 0 IBUF = 1 2 3 4 5 6 7 8 9 10
MYRANK = 1 IBUF = 1 2 3 0 5 6 7 0 0

```

# コミュニケータに関するサブルーチン MPI\_COMM\_SPLIT

## 機能 (関連する節: 3 - 6 節)

旧グループ(コミュニケータがcomm)に属する各プロセスを分割し、いくつかの新グループ(コミュニケータがnewcomm)を作成します。実行が終了すると、新グループに付けられた値がnewcommに戻ります。

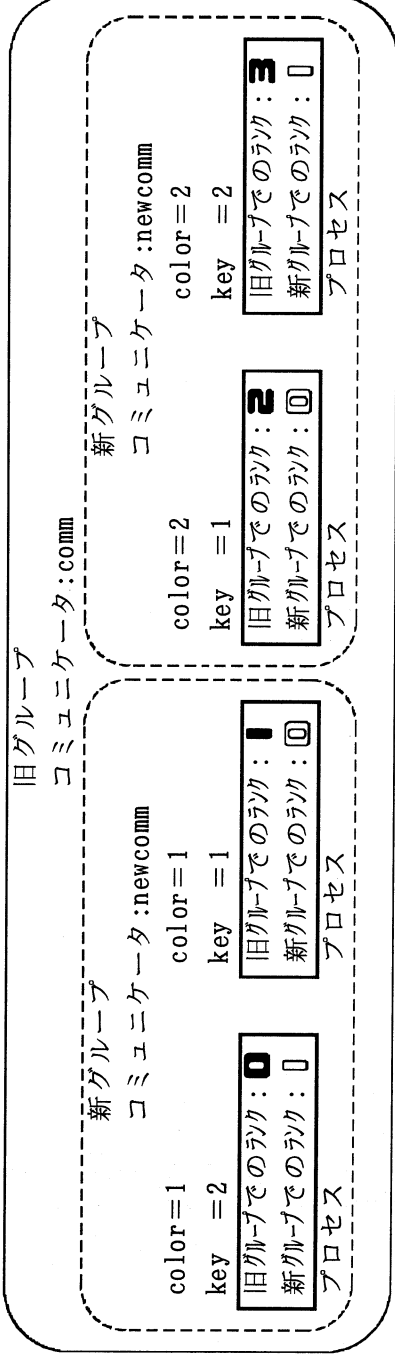
旧グループ(コミュニケータがcomm)に属する全プロセスが、本サブルーチンをコールする必要がある場合があります。  
colorで指定した値が同一であるプロセスが同じ新グループに所属し、新グループ内のランク(0,1,2,...)はkeyで指定した値の昇順に付けられます。

## 使用法

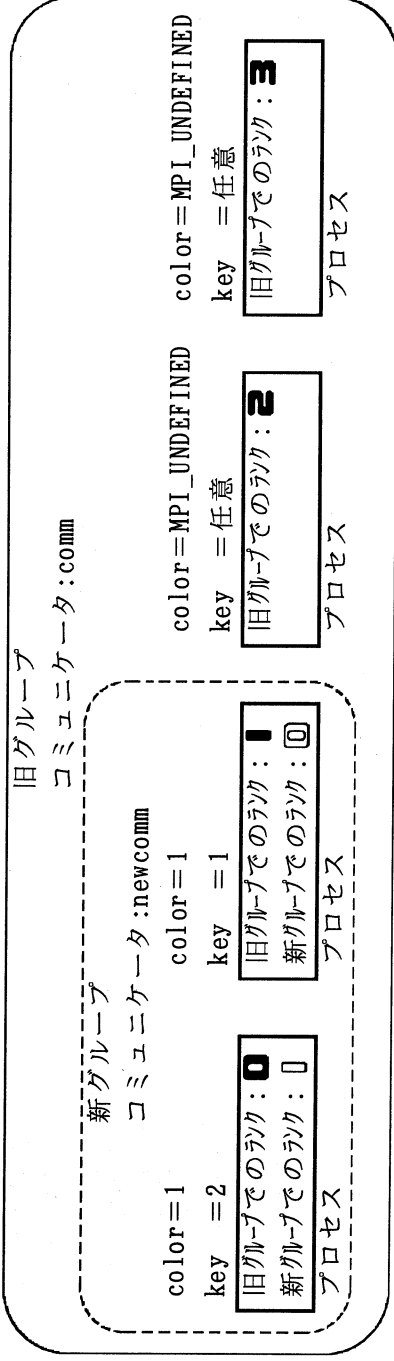
```
CALL MPI_COMM_SPLIT(comm, color, key, newcomm, ierror)
```

- comm : 整数。旧グループのコミュニケータ(グループ名)を指定します。
- color : 整数。ここで指定した値が同一であるプロセスが同じ新グループに所属します。ゼロ以上の整数値を指定して下さい。MPI\_UNDEFINEDを指定した場合には、そのプロセスは新グループに所属せず、newcommにはMPI\_COMM\_NULLが戻ります。
- key : 整数。ここで指定した値の昇順に、新グループ内のランクが0,1,2,...の順につけられます。keyの値が同一のプロセスが複数ある場合は、旧グループのランクの順にランクがつけられます。
- newcomm : 整数。作成する新グループのコミュニケータ(グループ名)を指定します。実行が終了すると、新グループに付けられた値が戻ります。
- ierror : 整数。完了コードが戻ります。正常終了の場合は「MPI\_SUCCESS」という値が戻ります。

### 【例1】



### 【例2】



## ■ サンプルプログラム 【例1】 のプログラム例です。

```

PROGRAM MAIN
INCLUDE 'mpif.h'
CALL MPI_INIT(IERR)
CALL MPI_COMM_SIZE(MPI_COMM_WORLD, NPROCS, IERR)
CALL MPI_COMM_RANK(MPI_COMM_WORLD, MYRANK, IERR)
IF (MYRANK == 0) THEN
  ICOLOR = 1
  IKEY = 2
ELSEIF (MYRANK == 1) THEN
  ICOLOR = 1
  IKEY = 1
ELSEIF (MYRANK == 2) THEN
  ICOLOR = 2
  IKEY = 1
ELSEIF (MYRANK == 3) THEN
  ICOLOR = 2
  IKEY = 2
ENDIF
CALL MPI_COMM_SPLIT(MPI_COMM_WORLD, ICOLOR, IKEY, NEW_WORLD, IERR)
CALL MPI_COMM_SIZE(NEW_WORLD, NEWPROCS, IERR)
CALL MPI_COMM_RANK(NEW_WORLD, NEWRANK, IERR)
PRINT *, 'MYRANK = ', MYRANK, 'NEWPROCS = ', NEWPROCS, 'NEWRANK = ', NEWRANK
CALL MPI_FINALIZE(IERR)
END

```

MPIが使用するインクルードファイルを指定します。  
MPI環境の初期化処理を行います。  
プロセス数NPROCSを取得します。  
自分のランクMYRANKを取得します。  
各プロセスは、新グループでのcolorとkeyの値を設定します。

新グループを作成します。  
新グループのプロセス数NEWPROCSを取得します。  
新グループでのランクNEWRANKを取得します。

MPI環境の終了処理を行います。

```

[aoyama@node01]~/aoyama: mpirun -np 4 a.out
MYRANK = 0 NEWPROCS = 2 NEWRANK = 1
MYRANK = 1 NEWPROCS = 2 NEWRANK = 0
MYRANK = 2 NEWPROCS = 2 NEWRANK = 0
MYRANK = 3 NEWPROCS = 2 NEWRANK = 1

```

## 参考文献

\*がついている資料は英語で書かれています。本書を改訂した影響で、文献の参照番号が不連続になっています。

- [ 1 ] 本書は <http://accr.riken.jp/HPC/training.html> に掲載されています。
- [ 2 ] Linuxで並列処理をしよう 第2版 石川裕 他 著 (共立出版)
- [ 3 ] 反復法の数理 藤野清次、張紹良 著 (朝倉書店)
- [ 4 ] 並列プログラミングの基礎 大森健児 訳 (丸善株式会社)
- [ 5 ] 並列計算法入門 樺山和男、西村直志、牛島省 著 (丸善)
- [ 6 ] チューニング技法入門 青山幸也 著 <http://accr.riken.jp/HPC/training/text.html>
- \*[15] MPI:A Message-Passing Interface Standard Message Passing Interface Forum MPI標準そのもので、<http://www.mpi-forum.org/> からダウンロードすることができます([27]に日本語訳があります)。MPI-2.0の第3章にMPI-1.1の正誤表があります。またこのサイトの「How to make comments」で、MPIに関する質問(英語)を受け付けています。
- [16] 超並列有限要素解析 矢川元基、塩谷隆二 著 (朝倉書店)
- [17] 計算力学とCAEシリーズ 7 パラレルコンピュータインテグ 矢川元基、曾根田直樹 著 (培風館) 有限要素法の領域分割法による並列化方法について説明されています。
- [18] COMシリーズ 超並列コンピュータ入門 野口正一 監修、村岡洋一、山名早人 著 (オーム社) ベクトル計算機と並列計算機の基本的な項目が説明されています。
- [19] 計算力学7ー計算力学における超並列計算法 矢川元基、奥田洋司 共編 (養賢堂)
- \*[20] USING MPI, Using MPI-2 (2分冊になりました)
- William Gropp, Ewing Lusk, Anthony Skjellum 著 (The MIT Press)
- [21] 実践MPI-2 ウィリアム・グロップ 他 著、畑崎隆雄 訳 (株式会社ピアソン・エデュケーション) [20]の第2巻の日本語訳です。
- \*[22] MPI: The Complete Reference (2分冊になりました) Marc Snir, Steve W.Otto, Steven Huss-Lederman, David W.Walker, Jack Dongarra 著 (The MIT Press)
- [23] 行列計算ソフトウェア WS,スーパーコン,並列計算機 小国力 編著、村田健郎、三好俊郎、ドンガラ,J.J.,長谷川秀彦 著 (丸善)
- [24] スーパーコンピュータ 科学技術計算への適用 村田健郎、小国力、唐木幸比古 著 (丸善)
- [27] <http://phase.hpcc.jp/phase/mpi-j/ml/> MPI仕様書[15]の日本語訳があります。MPI-2.0の第3章にMPI-1.1の正誤表があります。
- [28] <http://www.egroups.co.jp/group/mpi-j> MPIに関するメーリングリストです。
- [31] Twisted Data Layout 進藤達也、岩下英俊、土肥実久、萩原純一、金城シヨーン 著 「並列処理シミュレーションJSPP'94」ツイスト分割法についての論文です。
- [36] 入門Fortran 90 実践プログラミング 東田幸樹、山本芳人、熊沢友信 著 (SOFT BANK)
- \*[37] PARALLEL PROGRAMMING with MPI Peter S. Pacheco 著 (Morgan Kaufmann Publishers, Inc.) 例題は全てC言語で書かれています。Fortran版も <http://www.usfca.edu/mpi> にあります。
- [38] MPI並列プログラミング P.パチャエコ 著、秋葉博 訳 (培風館) [37]の日本語訳です。
- [39] 「超並列コンピュータの現状と動向調査」調査報告書 (富士総合研究所) 現在稼働中の並列コンピュータの特性や市場動向、並列アルゴリズムなどが紹介されています。
- [40] 富士総研技報 vol.5 No.1 1995 特集 超並列コンピュータインテグ (富士総合研究所) 現在稼働中の並列コンピュータの特性や市場動向、並列アルゴリズムなどが紹介されています。
- [41] MPIに関連するインターネットのリンクを以下に示します(参考文献[37]より転載)。
  - フリーウェアのMPI
  - <ftp://info.mcs.anl.gov> (mpich)
  - <ftp://ftp.osc.edu/pub/lam> (LAM)
  - <ftp://ftp.epcc.ed.ac.uk/pub/chimp/release> (CHIMP)
  - <ftp://csftp.unomaha.edu/pub/rewini/WinMPI> (WinMPI)
  - <http://pandora.ac.pt/w32mpi> (W32MPI)
  - <http://www.open-mpi.org/> (Open MPI)



- MPIのFAQ(frequently asked questions)です。  
<http://www.erc.msstate.edu/mpi/mpi-faq.html>
- MPI Web Page  
<http://www.mcs.anl.gov/mpi>  
<http://www.erc.msstate.edu/mpi>  
<http://www.epm.ornl.gov/~walker/mpi>

- ニュースグループ  
[comp.parallel.mpi](mailto:comp.parallel.mpi)  
[fj.comp.parallel](mailto:fj.comp.parallel)
- MPI-2, MPI-IO  
<http://www.mcs.anl.gov/Projects/mpi/mpi2/mpi2.html>  
<http://loveface.nas.nasa.gov/MPI-IO/mpi-io.html>

- 参考文献[37]のホームページとサンプルプログラム  
<http://www.mkp.com>

<http://www.usfca.edu/mpi>

#### [42] 学会/メーリングリスト/リンク集など

- PHASE(Parallel and HPC Application Software Exchange)  
<http://phase.hpcc.jp/>

HPC関係の各種リンクやメーリングリストのリンクがあります。

- SWoPPメーリングリスト  
[http://www.hpcc.jp/swopp/ml\\_readme.html](http://www.hpcc.jp/swopp/ml_readme.html)
- SofTek HPCメーリングリスト  
<http://www.softek.co.jp/>
- SMPP(超並列計算研究会)  
<http://www.is.doshiha.ac.jp/SMPP/>

「PCクラスタ超入門」というダウンロードできる資料が置いてあります。

- CFDネット(数値流体力学関係のメーリングリスト)

[majordomo@cf.d.ritsumei.ac.jp](mailto:majordomo@cf.d.ritsumei.ac.jp)宛てに『subscribe cfdnet 登録したいメールアドレス』だけを書いたメールを送付すると、管理者が後で登録します。

- HPC関係のニュースを集めたサイト

<http://ytaka.at.infoseek.co.jp/>

- 日本計算工学会

<http://www2.kajima.co.jp/jscs/wwwjscs/jscs.html>

- 情報処理学会HPC研究会

<http://phase.etl.go.jp/sighpc/>

- 新情報処理開発機構(並列処理シンポジウムへのリンクあり)

<http://www.rwcp.or.jp/>

- 日本原子力研究開発機構 システム計算科学センター

<http://ccse.jaea.go.jp/ja/index.html>

- SWoPP(並列/分散/協調処理に関するサマリー・ワークショップ)

<http://www.hpcc.jp/swopp/>

#### [43] 数値計算ライブラリー

- PCP(離散化並列解法のための並列計算プラットフォーム)

<http://www.ibase.aist.go.jp/infobase/pcp/index.html>

- SSI(大規模シミュレーション向け基盤ソフトウェア)

<http://www.ssisc.org/>

- ScaLAPACK

<http://www.netlib.org/scalapack/index.html>

以下に、LAPACKとScaLAPACKに関する日本語のドキュメントがあります。

<http://phase.hpcc.jp/phase/lapack-j/index.html>

<http://phase.hpcc.jp/phase/lapack-j/ScaLAPACK/index.htm>

<http://spring.cc.kyushu-u.ac.jp/scp/system/manual/ScaLAPACK/scalapack.pdf>

- PETSc <http://www-unix.mcs.anl.gov/petsc/petsc-2/>
- Lis <http://www.ssisc.org/lis/>
- PARCEL <http://guide.tokai.jaeri.go.jp/program/software/list/pg1/index.html>
- GeoFEM <http://geofem.tokyo.rist.or.jp/>
- FFTW <http://www.fftw.org/> (高速フーリエ変換を高速に行うFFTWのサイト)
- ICCG法の並列版(京都大学大型計算機センター)
  - <http://fem.kuee.kyoto-u.ac.jp/~take/parasolver/piccg/index.html> (〓は波線です)
- NAG <http://www.nag-j.co.jp/> (日本ニューメリカルゴリズムズグループ(株))
- IMSL <http://www.vni-j.co.jp/> (日本ビジュアルニューメリックス(株))
- SMS(行列計算ライブラリー) (MPI版は開発中?)
  - <http://www.vinas.com/seihin/sms/index.html> ((株)ヴァイナス)
- \*[44] On parallelism and convergence of incomplete LU factorizations, Shun Doi, Applied Numerical Mathematics 7(1991) 417-436, North-Holland
- \*[45] Introduction to Parallel Computing, Design and Analysis of Algorithms, Vipin Kumar, Ananth Grama, Anshul Gupta, George Karypis, The Benjamin/Cummings Publishing Company, Inc.
- [46] bit別冊 はじめての並列プログラミング 湯淺太一、安村通晃、中田登志之 編 (共立出版)
- [47] たのしくできる並列処理コンピュータ 小畑正貴 著 (東京電機大学出版局)
- [48] つくる並列処理コンピュータ 小畑正貴 著 (東京電機大学出版局)
- \*[49] SOLVING PARALLEL DIFFERENTIAL EQUATIONS ON PARALLEL COMPUTERS  
JIANGING ZHU 著 (World Scientific) 差分法の並列アルゴリズムに関する良書です。
- [50] NUMERICAL RECIPES in C [日本語版] (技術評論社)
- [51] ICCG法(SCG法)の説明 <http://www.ppc.nec.co.jp/> で、左欄の「登録ソフト」を選択し、ソフト名「psc98」を選択して下さい。
- \*[52] PARALLEL PROGRAMMING BARRY WILKINSON, MICHAEL ALLEN 著 (PRENTICE HALL)
- [53] 並列プログラミング入門 飯塚肇、緑川博子 共訳 (丸善)
  - [52]の第I部の日本語訳です。
- [54] <http://guide.tokai.jaeri.go.jp/program/> 日本原子力研究所 計算科学技術推進センターが公開している並列関係の資料(論文や並列ライブラリーなど)です。
- [55] 並列有限要素解析[I] クラスタコンピュータライティング 奥田洋司、中島研吾 共編 (培風館)
- [56] 並列有限要素解析[II] 並列構造解析ソフトウェアFrontSTRを使いこなす 奥田洋司編著 (培風館)
- \*[57] Parallel Programming in OpenMP Rohit Chandra 他 著 (MORGAN KAUFMANN PUBLISHERS)
- \*[58] Parallel Programming in C with MPI and OpenMP Michael J. Quinn 著 (McGraw-Hill)
- \*[59] High Performance LINUX CLUSTERS with OSCAR, Rocks, openMosix, & MPI Joseph D. Sloan 著
- [60] OpenMPの仕様書(英語) <http://openmp.org/wp/openmp-specifications/>
- [61] OpenMPの仕様書([60]の日本語訳) [60]で「Version 3.0,Japanese」を選択して下さい。
- [62] OpenMPによる並列プログラミングと数値計算法 牛島省 著 (丸善)
- [63] OpenMP入門 南里豪志 著 <http://www.cc.kyushu-u.ac.jp/scp/system/library/OpenMP/OpenMP.html>
- [64] OpenMP入門(1)~(4) 計算工学 (2006No.4, 2007No.1,2,3) 南里豪志 著 (日本計算工学会)
- [65] マルチコアCPUのための並列プログラミング 並列処理&マルチスレッド入門  
安田絹子、小林林広、飯塚博道、阿部貴之、青柳信吾 著 (秀和システム)
- [66] C/C++プログラマーのためのOpenMP並列プログラミング 菅原清文 著 (カトシステム)
- [67] OpenMP入門 北山洋幸 著 (秀和システム)
- \*[68] Using OpenMP: Portable Shared Memory Parallel Programming Barbara Chapman 他 著 (The MIT Press)
- [69] 「並行コンピュータライティング技法」 千住 治郎 訳 (オライリー・ジャパン)
- [70] 「並列数値処理 - 高速化と性能向上のために -」 金田 康正 編著 (コロナ社)
- [71] 「コンピュータ設計の基礎」 Hisa Ando 著 (マイコミジャーナルブックス)
- [72] 「プロセッサを支える技術」 Hisa Ando 著 (技術評論社)
- [73] 「高性能コンピュータ技術の基礎」 Hisa Ando 著 (マイコミジャーナルブックス)