

計算科学演習B MPI 発展

学術情報メディアセンター 情報学研究科・システム科学専攻 中島 浩

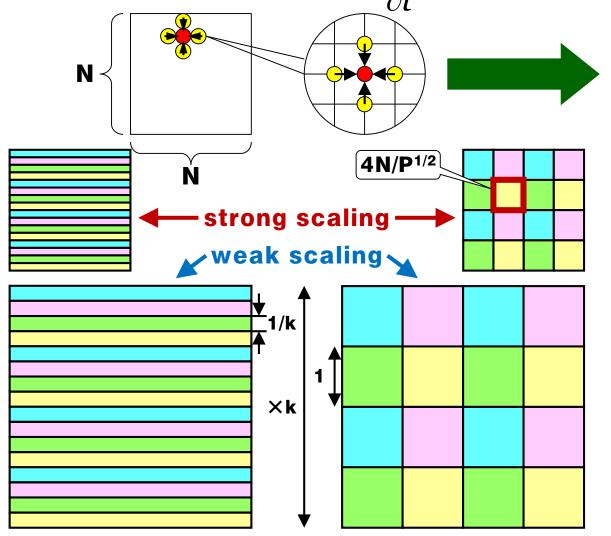


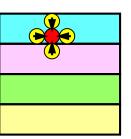
- 序論
 - 1次元分割 vs 2次元分割、実装のポイント
- プロセスの2次元配置
 - 直交プロセス空間
 - rank ⇔ 座標の変換
- 2次元配列の宣言・割付・参照
- 不連続データの通信
 - 派生データ型、東西通信用派生データ型
 - 通信回数削減との関係
- 2次元分散データの出力
 - 良くない方法2種
 - 個別出力:MPI File I/O 機能の利用

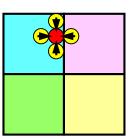


(注) 1次元分割 vs 2次元分割

拡散方程式 $\nabla^2 \varphi = \frac{\partial \varphi}{\partial t}$ の初期値問題求解 by 陽解法





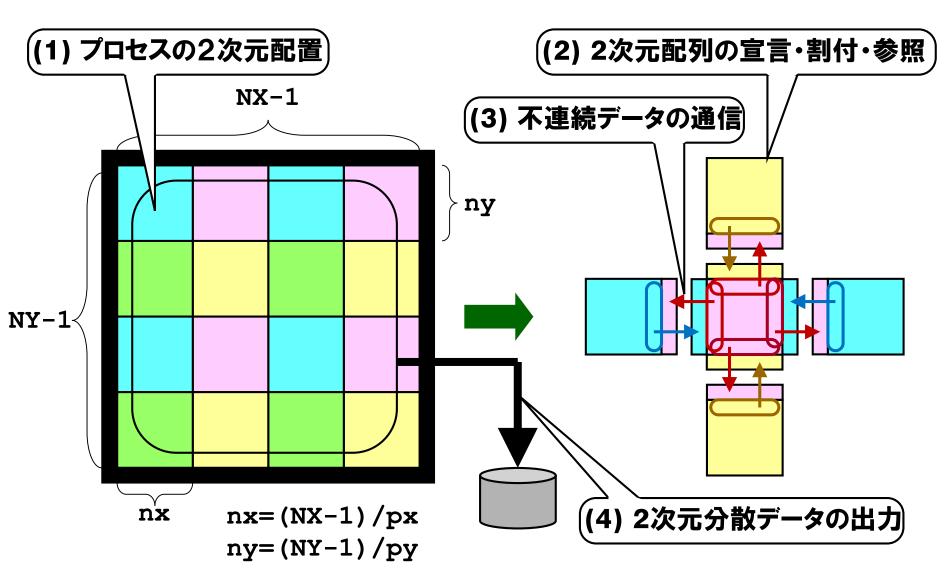


1次元分割

2次元分割

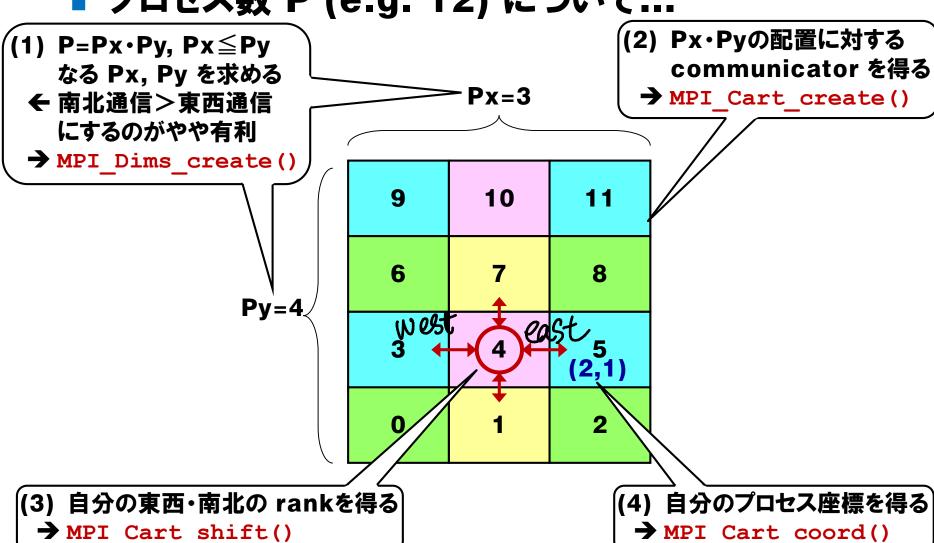
どちらも ...

- 部分問題計算量 ≒1/P
- 部分問題メモリ量 ≒ 1 / P (留意点後述)
- 分割不能計算量 ≒0 (留意点後述)
- 計算/step=O(N²/P) 通信/step
 - $= O(N) \text{ or } O(N/P^{1/2})$
 - → 通信/計算<0(1)</p>



プロセスの2次元配置機論

■ プロセス数 P (e.g. 12) について...



プロセスの2次元配置 直交プロセス空間 (1/2)

- n=n₀×n₁×...×n_{k-1} (n₀≥n₁≥...≥n_{k-1} かつ Σn_i 最小) の算出
 - MPI_Dims_create(int n, int k, int *dims)
 dims={0,...,0} or dims=(/0,...,0/) で呼出し
 → dims[0...k-1]=dims(1...k)=n₀, n₁, ..., n_{k-1}
 cf) 呼出し前の非ゼロ要素は「尊重」される
 - ■例

```
int dims[2]={0,0};
MPI_Comm_size(MCW,&np);
MPI_Dims_create(np,2,dims);
```

```
integer::dims(2) = (/0,0/)
call MPI_Comm_size(MCW,np,err)
call MPI_Dims_create(np,2,dims,err)
```

プロセスの2次元配置 直交プロセス空間 (2/2)

- プロセス座標系(p₀, ..., p_{k-1}) (0≦p_i<n_i) 生成
 - - $dims=n_0, n_1, ..., n_{k-1}$

 - reorder==0 ? rank不変 : rank変更あり

- 例

```
int periods[2]={0,0}; MPI_Comm cart;
MPI_Dims_create(np,2,dims);
MPI_Cart_create(MCW,2,dims,periods,0,&cart);
```

```
integer::periods(2)=(/0,0/),cart
call MPI_Dims_create(np,2,dims,err)
call MPI_Cart_create(MCW,2,dims,periods,0,cart,err)
```

プロセスの2次元配置 rank ⇔ 座標の変換 (1/3)

- (p₀, ..., p_{k-1}) Ø rank (row major)
 - $\mathbf{r}_0 = \mathbf{p}_0 \quad \mathbf{r}_i = \mathbf{n}_i \mathbf{r}_{i-1} + \mathbf{p}_i \quad \mathbf{r}_{k-1} = \mathbf{rank}(\mathbf{p}_0, ..., \mathbf{p}_{k-1})$

 $coord=p_0, ..., p_{k-1} \rightarrow rank=rank(p_0, ..., p_{k-1})$

- 例

```
int c[2];
for(c[0]=0;c[0]<dims[0];c[0]++)
  for(c[1]=0;c[1]<dims[1];c[1]++) {
     MPI_Cart_rank(cart,c,&r); MPI_Recv(...,r,...);
  }
integer::c(2)
do c(1)=0,dims(1)-1; do c(2)=0,dims(2)-1
  call MPI_Cart_rank(cart,c,r,err)
  call MPI_Recv(...,r,...)
end do; end do</pre>
```

プロセスの2次元配置 rank ⇔ 座標の変換 (2/3)

- 逆変換 (p₀, ..., p_{k-1})=rank⁻¹(r)
 - MPI_Cart_coords(MPI_Comm cart, int rank, int k, int *coord)

- 例

```
int c[2];
MPI Cart coords(cart,me,2,c);
for(i=...) for(j=...) {
  v=c[0]*nv+i+1; x=c[1]*nx+j+1;
}
integer::c(2)
call MPI Cart coords (cart, me, 2, c, err)
do i=\ldots; do j=\ldots;
  y=c(1)*ny+i+1; x=c(2)*nx+j+1
end do; end do
```

プロセスの2次元配置 rank ⇔ 座標の変換 (3/3)

■ 隣接プロセスの rank

```
MPI Cart shift (MPI Comm cart, int d, int dir,
                    int *src, int *dst)
 me=rank(p_0, ..., p_{k-1}) \rightarrow src=rank(p_0, ..., p_d \mp 1, ..., p_{k-1})
                         dst=rank(p<sub>0</sub>, ..., p<sub>d</sub>±1, ..., p<sub>k-1</sub>)
       非周期境界の端では
       MPI PROC NULL
MPI Cart shift(cart,0,1,&south,&north);
MPI Cart shift(cart,1,1,&west,&east); ...
MPI Sendrecv(...,north,...,south,...); ...
MPI Sendrecv(...,east,...,west,...); ...
 call MPI Cart shift(cart, 0, 1, south, north, err)
 call MPI Cart shift(cart,1,1,west,east,err); ...
 call MPI Sendrecv(..., north,..., south,...); ...
 call MPI Sendrecv(..., east,..., west,...); ...
```



2次元配列の宣言・割付・参照 (1/3)

- [X0,x1]×[Y0,y1] なる配列 a の宣言&割付 (X0, Y0 は定数、x1, y1は可変)
 - Fortran は相変わらず簡単

```
double precision,allocatable::a(:,:)
allocate(a(X0:x1,Y0:y1))
```

ANSI C で a[j][i] に固執→良くない方法しかない

```
double **a; int h=y1-Y0+1, w=x1-X0+1;
a=(double**)malloc(h*sizeof(double*))-Y0;
a[Y0]=(double*)malloc(h*w*sizeof(double))-X0;
for(j=Y0+1;j<=y1;j++) a[j]=a[j-1]+w;</pre>
```



2次元配列の宣言・割付・参照 (2/3)

■ ANSI Cでの良い方法=2次元配列の1次元化

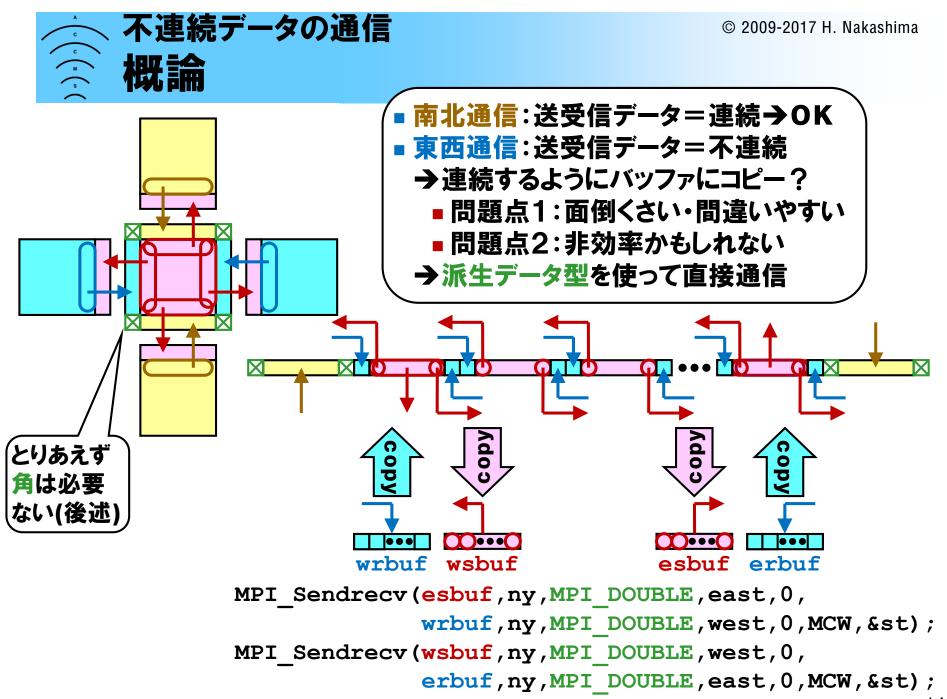
```
double *a;
 a=(double*)malloc(h*w*sizeof(double))-Y0*w-X0;
 for (i=...) for (j=...)
    ... = ... a[i*w+j]... a[(i+1)*w+j]...;
→内側ループのループ不変式 (e.g. i*w) の追い出し
 for(i=...) {
   double *ai=a+i*w, *aip=ai+(i+1)*w;
   for (j=...)
      ...=...ai[j]...aip[j]...;
→演算子強度低減 (strength reduction) 乗算→加算
 for (i=..., ai=..., aip=...; ...; ..., ai+=w, aip+=w) {
   for(j=...)
      ...=...ai[j]...aip[j]...;
```



2次元配列の宣言・割付・参照 (3/3)

C99での良い方法=可変長配列(!!)

```
ローカル変数または
• f(...) {...
                                 ローカル変数
   h=...; w=...; |引数のみ可能
                                 →h*w が大きい
                                 →スタックオーバフロー
   double b[h][w]
   double (*a)[w] = (double(*)[w]) &b[-Y0][-X0];
 1次元目が ₩ 要素の
                    引数の場合
 配列へのポインタ
                    f(int h, int w, double a[h][w]){
                     →要素数を定める引数よりも後で定義
• f(...) {...
   h=\ldots; w=\ldots;
   double (*a)[w];
   a=(double(*)[w])malloc(h*w*sizeof(double));
   a=(double(*)[w])&a[-Y0][-X0]
                                |配列の実体はヒープに
  . . . }
                                →h*w 大でもOK
■ コンパイルは mpiicc -std=c99 ...
```





- 派生データ型 (derivative data type)
 - 基本型 (MPI_DOUBLE など)・他の派生型を組み合わせて 作られる送受信用 (およびファイルI/O用) の型
 - 内部に穴がある不連続データ型も定義可能
 - →送受信MPI関数が必要に応じて pack/unpack
 - 規則的な派生データ型生成関数

東西通信

- MPI_Type_contiguous() 同一型の単純な連続
- MPI Type vector() 同一型の規則的な連続+不連続
- MPI_Type_hvector() 同一型の規則的な連続+不連続
- 不規則な派生データ型生成関数

結果出力

- MPI Type_indexed() 同一型の不規則な連続+不連続
- MPI_Type_hindexed() 同一型の不規則な連続+不連続
- MPI_Type_create_subarray()

同一型の不規則な連続+不連続

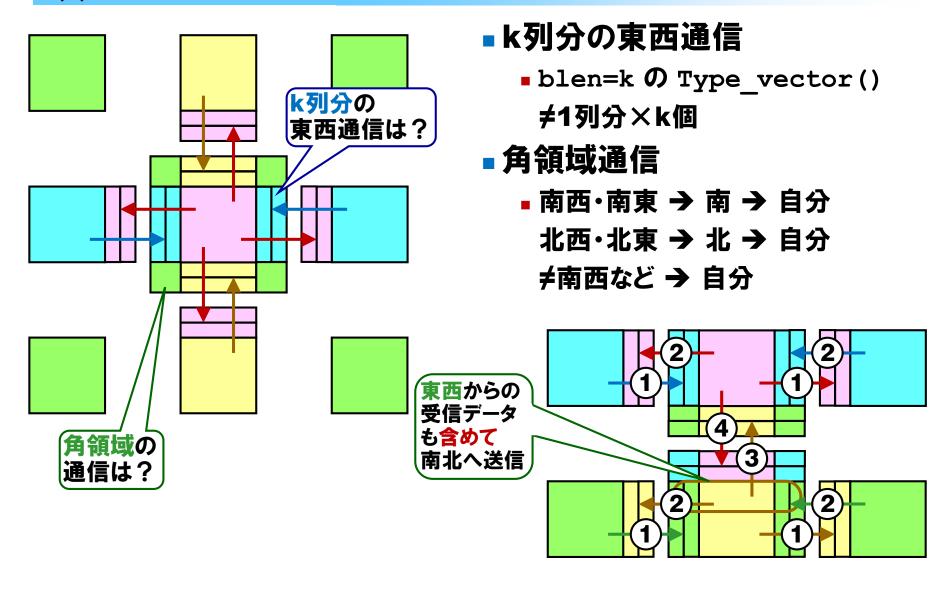
■ MPI Type struct() 非同一型の不規則配置

不連続データの通信

東西通信用派生データ型

MPI Type vector(int count, int blen, int stride, MPI Datatype oldtype, MPI Datatype *newtype) blen個 oldtype stride個分 $\bullet \bullet \bullet$ count個 MPI DOUBLE ×1個 nx+2ny個 MPI Datatype vedge; MPI Type vector(ny,1,nx+2,MPI DOUBLE,&vedge); MPI Type commit(&vedge); ____ 作ったデータの使用開始宣 MPI_Sendrecv(&u[...],1,vedge,east,0, &u[...],1,vedge,west,0,MCW,&st); integer::vedge; call MPI Type vector(ny,1,nx+2,MPI DOUBLE PRECISION, vedge,err); call MPI Type commit(vedge,err); call MPI Sendrecv(u(...),1,vedge,east,0, u(...),1,vedge,west,0,MCW,st,err);

不連続データの通信 通信回数削減への対応



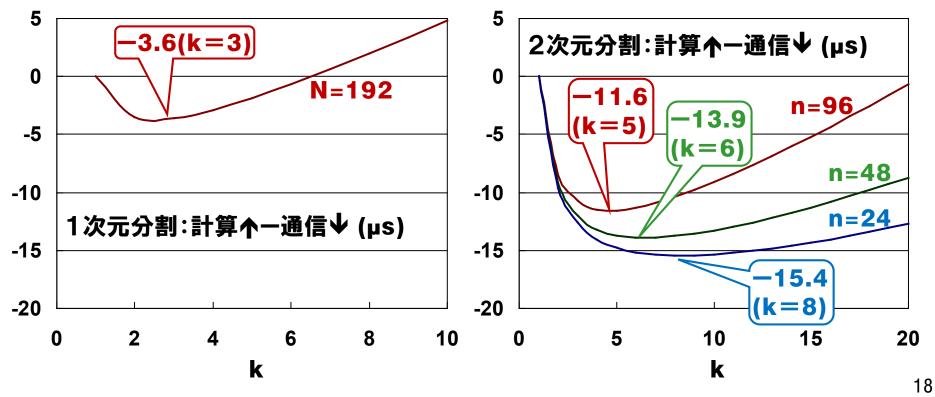


不連続データの通信

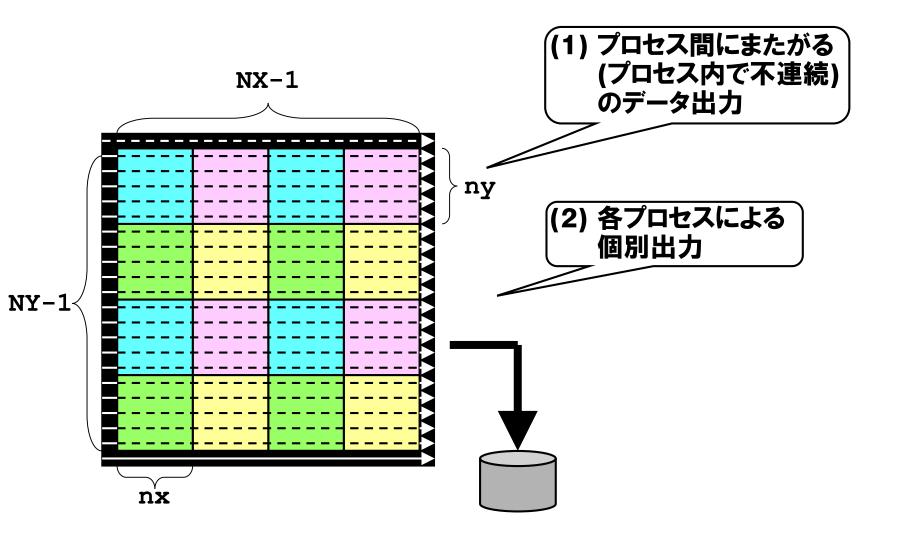
不理続アーダの通信通信回数削減の効果

分割	1次元	2次元
計算个	k(k-1)Nt/k	2k(k-1)(n+(2k-1)/3)t/k
通信Ψ	2(k-1)L/k	(4(k-1)L-4(k-1)(k-1+n)/B)/k



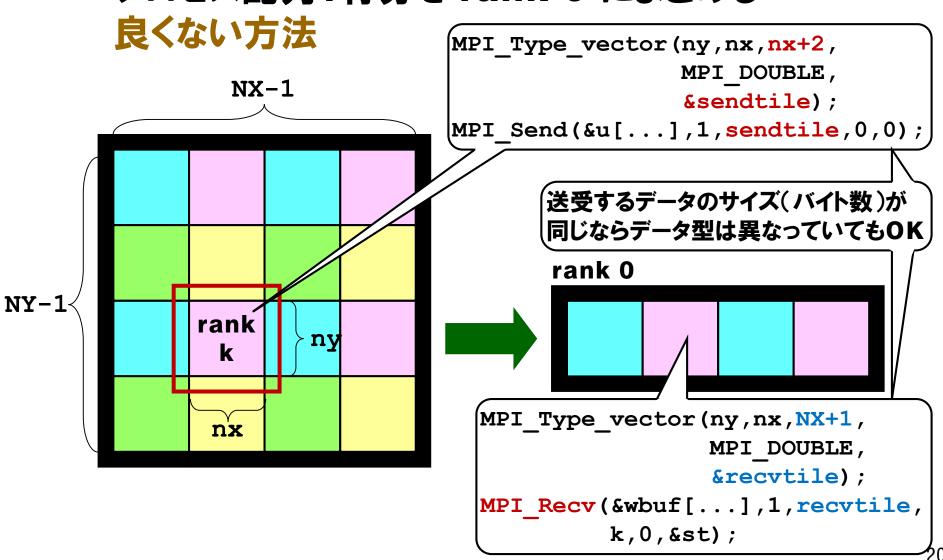


● 2次元分散データの出力 概論



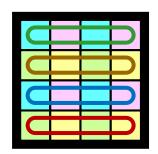
2次元分散データの出力 rank0 による出力

■ プロセス配列1行分を rank 0 にまとめる



2次元分散データの出力 西端プロセスによる出力

- プロセス配列1行分を西端にまとめる良くない方法
 - 同じ行のプロセスからなる communicator 生成

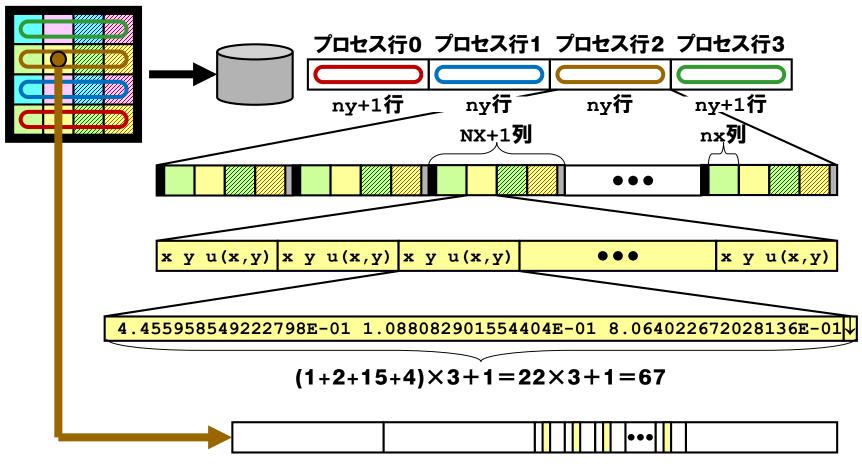


- newcomm = color が自分と同じプロセスの集合
- newcomm での rank=key の昇順
- 行プロセス集合の中で西端に MPI_Gather()
- 例:

```
MPI_Comm row;
MPI_cart_coords(cart,me,2,c);
MPI_Comm_split(cart,c[0],c[1],&row);
MPI_Gather(&u[...],...,wbuf,...,0,row);
```

2次元分散データの出力 個別出力:概論

やりたいこと=プロセスの位置により定まるファイル 中の不連続領域への分散書込み



個別出力:ファイルの open/close (1/2)

- open等のMPI関数:全プロセスが同一引数で実行
 - MPI File open (MPI Comm comm, char *name, int mode, MPI Info info, MPI File *fh) ファイルに関するヒント
 - comm で共有する名前が name のファイルをオープン
 - アクセスモード mode = 以下 (など) の定数の論理和

MPI MODE RDONLY: 読出し専用

MPI MODE RDWR :読出し&書込み

MPI MODE WRONLY:書込み専用

MPI MODE APPEND:追加(書込み)

MPI MODE CREATE:ファイルがなければ生成

- MPI File set size(MPI File fh, MPI Offset size)
 - ファイルの (初期) サイズを設定
 - 書込みオープンでも既存ファイルの初期サイズは0ではない!!
- MPI File close (MPI File *fh)

|書込み&存在しなければ生成

なった場合にゴミが残る

(全) 個別出力:ファイルの open/close (2/2)

Cの例

```
MPI File udata;
 MPI File open(cart,"...",
               MPI MODE WRONLY | MPI MODE CREATE,
特にヒントはなし
               MPI INFO NULL, &udata);
 MPI File set size(udata, (MPI Offset) 0);
 MPI_File_close(&udata);
                          存在した場合初期サイズを0に
                          しないと元より小さいファイルに
```

■ Fortran の例

integer::udata

```
integer(kind=MPI OFFSET KIND)::zero=0
call MPI File open(cart, '...', &
          ior (MPI MODE WRONLY, MPI MODE CREATE), &
          MPI INFO NULL, udata, err)
call MPI File set size(udata, zero, err)
call MPI File close (udata, err)
```

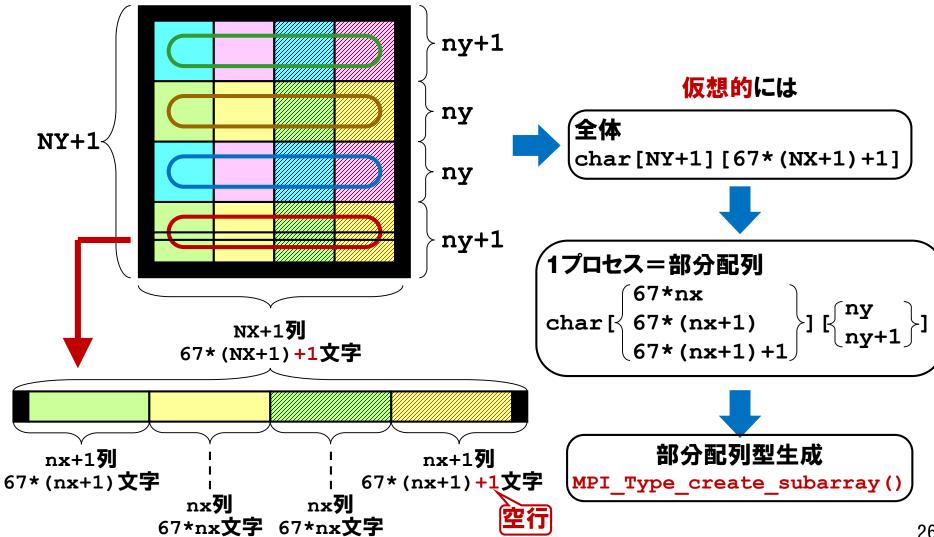
個別出力:個別書込みパターンの設定

■ ファイルの view の設定

MPI File set view (MPI File fh, MPI Offset disp, MPI Datatype etype, MPI Datatype ftype, char *drep, MPI Info info) ファイル形式情報 ファイルアクセスの 最小単位の型 京大スパコンでは etype "native" プロセス固有の P₀ アクセスパターン を示すデータ型 √(filetype) ftypefh #2 #3 #4 #5 #6 disp 1つの ftype を超えると ftype の繰り返し

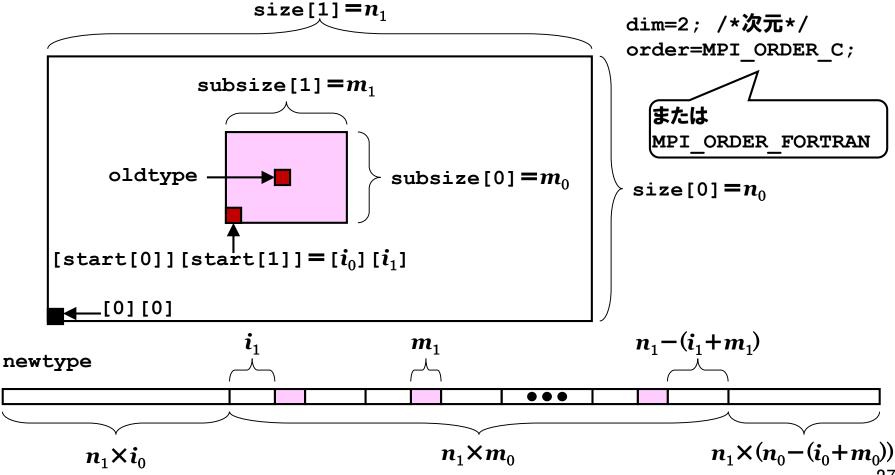
個別出力: filetype の生成 (1/4)

■ 出力ファイル=仮想的2次元配列



個別出力: filetype の生成 (2/4)

MPI Type create subarray(int dim, int *size, int *subsize, int *start, int order, MPI Datatype oldtype, MPI Datatype *newtype)





両方とも

成立する

両方とも

成立する

© 2次元分散データの出力 © 2 個別出力: filetype の生成 (3/4)

Cの例

```
#define LW 67
         MPI Datatype ftype;
         int size[2]={NY+1,LW*(NX+1)+1}, subsize[2], start[2];
         MPI Dims create(np,2,dims);
         MPI Cart coords (cart, me, 2, c);
         subsize[0]=ny; subsize[1]=LW*nx;
         start[0]=c[0]*ny+1; start[1]=LW*(c[1]*nx+1);
         if (c[0]==0) { subsize[0]++; start[0]=0; }
ことがある
         if/ (c[0] == dims[0]-1) subsize[0]++;
         if/ (c[1] == dims[1] -1) subsize[1] += LW+1;
ことがある
         MPI Type create subarray(2, size, subsize, start,
                                 MPI ORDER C,MPI CHAR, &ftype);
         MPI Type commit(&ftype);
         MPI File set view(udata, (MPI Offset) 0, MPI CHAR, ftype,
                           "native", MPI INFO NULL);
```

両方とも

成立する

両方とも

成立する

② 2次元分散データの出力 © 2 個別出力:filetype の生成 (4/4)

■ Fortran の例

```
integer,parameter::LW=67
           integer::ftype
           integer::size(2) = (/NNY+1, LW*(NNX+1)+1/),
           integer::subsize(2), start(2)
           call MPI Dims create(np,2,dims,err)
           call MPI Cart coords (cart, me, 2, c, err)
           subsize(1) = ny; subsize(2) = LW*nx;
           start(1) = c(1) * ny + 1; start(2) = LW*(c(2) * nx + 1)
          (c(1)==0) then
             subsize(1) = subsize(1) + 1; start(1) = 0; endif
ことがある
              (c(1) == dims(1) -1) subsize(1) = subsize(1) +1
             (c(2) == 0) then
             subsize(2) = subsize(2) + LW; start(2) = 0; endif
ことがある
           if/(c(2) == dims(2) - 1) subsize(2) = subsize(2) + LW+1
           call MPI Type create subarray(2, size, subsize, start, &
                             MPI_ORDER C,MPI CHARACTER,ftype,err)
           call MPI Type commit(ftype,err)
           call MPI File set view(udata, zero, MPI CHARACTER, &
                             ftype, 'native', MPI INFO NULL, err)
```

MPI File write()を

2次元分散データの出力

(1/3) **個別出力:ファイルの書込み** (1/3)

- MPI File write (MPI File fh, void *buf, int count, MPI Datatype dtype, MPI Status *st)
 - buf から count 個の dtype データを fh に書込み
 - dtype のデータがある部分を繋いだもの = etype を並べたもの
- あまり良くない (かもしれない) 例

```
沢山実行すると、とりあえず
char wbuf[LW+1];
                                 京大スパコンでは結構な時間
                                 がかかった
for(i=...) {
  for(j=...) {
    sprintf(wbuf,"...,u[...]);
    MPI File write(udata, wbuf, LW, MPI CHAR, &st);
  if(c[1] == dims[1] - 1)
    MPI File write(udata,"\forman n",1,MPI CHAR,&st);
```

(2/3) 個別出力:ファイルの書込み (2/3)

どのマシンでもそこそこ良い (はずの) 例 (C版)

```
char *wbuf;
wbuf=(char*)malloc((LW*(nx+2)+2)*sizeof(char));
for(i=...) {
  for (j=...,k=0;...,k+=LW)
    sprintf(wbuf+k,"...,u[...,u[...]);
  if(c[1] == dims[1] - 1)
    sprintf(wbuf+(k++),"\frac{\text{Y}}{n}");
  MPI File write(udata, wbuf, k, MPI CHAR, &st);
```

(3/3) 個別出力:ファイルの書込み(3/3)

どのマシンでもそこそこ良い (はずの) 例 (Fortran版)

```
subroutine print u(u,...,bufsize,eastmost)
integer::bufsize
logical::eastmost
character(len=bufsize)::wbuf
. . .
do i=\ldots; k=1
  do j=...
    write(wbuf(k:k+LW-1),'(...)')...,u(...),char(10)
    k=k+LW
  end do
  if (eastmost) then
    wbuf (k:k) = char (10)
    call MPI File write (udata, wbuf, k, MPI CHAR, st, err)
  else
    call MPI File write (udata, wbuf, k-1, MPI CHAR, st, err)
end do
end subroutine
```

2次元分散データの出力 個別出力:まとめ (1/2)

■ Cのまとめ

```
MPI File udata;
MPI Datatype ftype;
MPI File open(&udata,...);
MPI File set size(udata, (MPI Offset) 0);
/* ftype の生成 */
MPI Type commit(&ftype);
MPI File set view(udata,...,ftype,...);
for(i=...) {
  for(j=...) { ... } ...;
  MPI File write(udata,...);
MPI File close(&udata);
```

Fortran のまとめ

```
integer::udata,ftype
...

call MPI_File_open(udata,...)

call MPI_File_set_size(udata,zero,err)

/* ftype の生成 */

call MPI_Type_commit(ftype,err)

call MPI_File_set_view(udata,...,ftype,...)

call print_u(u,...,LW*(nx+2)+1,c(1)==dims(1)-1)

call MPI_File_close(udata,err)
```