

OpenMP 並列化

- L2-solをOpenMPによって並列化する。
 - 並列化にあたってはスレッド数を「PEmpTOT」によってプログラム内で調節できる方法を適用する
- 基本方針
 - 同じ「色」(または「レベル」)内の要素は互いに独立, したがって並列計算(同時処理)が可能

科学技術計算のための マルチコアプログラミング入門 第三部: OpenMPによる並列化

2011年6月30日・7月1日

中島研吾

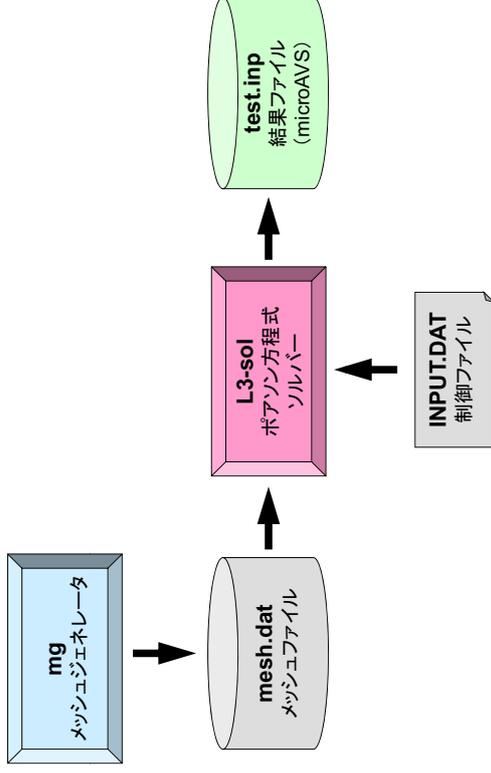
プログラムのありか

- 所在
 - <L3>/src, <L3>/run
- コンパイル, 実行方法
 - 本体
 - cd <L3>/src
 - Make
 - <L3>/run/l3-sol (実行形式)
- メッシュ生成
 - cd <L3>/run
 - f90 -O mg.f, cc -O mg.c
- コントロールデータ
 - <L3>/run/INPUT.DAT
- 実行用シェル
 - <L3>/run/x0.sh, x5.sh, x6.sh

実行例

```
% cd <L3>
% ls
  run  src  ft  reorder  hid
% cd src
% make
% cd ../run
% ls L3-sol
  L3-sol
% f90 -O3 mg.f -o mg
% <modify "INPUT.DAT">
% <modify "x0.sh">
% qsub x0.sh
```

プログラムの実行 プログラム, 必要なファイル等



制御データ(INPUT.DAT)

```

1.00e-00 1.00e-00 1.00e-00 DX/DY/DZ
1.0e-08 EPSICCG
16 PEsmptOT
100 NCOLortot
  
```

変数名	型	内容
DX, DY, DZ	倍精度実数	各要素の3辺の長さ (ΔX, ΔY, ΔZ)
EPSICCG	倍精度実数	収束判定値
PEsmptOT	整数	データ分割数
NCOLortot	整数	Ordering手法と色数 ≥2 : MC法 (multicolor) , 色数 =0 : CM法 (Cuthill-McKee) =-1 : RCM法 (Reverse Cuthill-McKee) ≤-2 : CM-RCM法

- L2-solへのOpenMPの実装
- Hitachi SR11000/J2での実行
 - 計算結果
 - CM-RCMオーダリング
- T2K(東大)での実行
- T2K(東大)での性能向上への道
 - NUMA Control
 - First Touch
 - データ再配置

L2-solにOpenMPを適用

- ICCGソルバーへの適用を考慮すると
- 内積, DAXPY, 行列ベクトル積
 - もともとデータ依存性無し ⇒ straightforwardな適用可能
- 前処理(修正不完全コレスキー分解, 前進後退代入)
 - 同じ色内は依存性無し ⇒ 色内では並列化可能

実はこのようにしてDirectiveを
直接挿入しても良いのだが・・・(1/2)

```
!$omp parallel do private(i, VAL, k)
do i = 1, N
  VAL= D(i)*W(i, P)
do k= indexL(i-1)+1, indexL(i)
  VAL= VAL + AL(k)*W(itemL(k), P)
enddo
do k= indexU(i-1)+1, indexU(i)
  VAL= VAL + AU(k)*W(itemU(k), P)
enddo
W(i, Q)= VAL
enddo
!$omp end parallel do
```

- スレッド数をプログラムで制御できるようにしてみよう

実はこのようにしてDirectiveを
直接挿入しても良いのだが・・・(2/2)

```
do icol= 1, NCOLORTot
!$omp parallel do private(i, VAL, k)
do i= COLORindex(icol-1)+1, COLORindex(icol)
  VAL= D(i)
do k= indexL(i-1)+1, indexL(i)
  VAL= VAL - (AL(k)**2) * DD(itemL(k))
enddo
DD(i)= 1. d0/VAL
enddo
!$omp end parallel do
enddo
```

- スレッド数をプログラムで制御できるようにしてみよう

ICCG法の並列化: OpenMP

- 内積: OK
- DAXPY: OK
- 行列ベクトル積: OK
- 前処理

プログラムの構成(1/2)

```
program MAIN
use STRUCT
use PGG
use solver_ICCG.mc
implicit REAL*8 (A-H, O-Z)
real (kind=8), dimension(:), allocatable :: WK
call INPUT
call POINTER_INIT
call BOUNDARY_CELL
call CELL_METRICS
call POI_GEN
PHI= 0. d0
call solve_ICCG.mc
& ( CELTOT, PHIL, AL, NCOLORTot, itemU, itemL, D,
& BEFORCE, PHIL, AL, NCOLORTot, PcompTOT,
& SMPindex, SMPindexG, EPSICCG, ITR, IER)
```

プログラムの構成 (2/2)

```
allocate (WK (ICELTOT))
do ioc= 1, ICELTOT
  ioc1= NEItoOLD(ioc)
  WK(ioc1)= PHI(ioc)
enddo
do ioc1= 1, ICELTOT
  PHI(ioc1)= WK(ioc1)
enddo
call OUTUCD
stop
end
```

結果 (PHI) をもとの番号
付けに もとす

プログラムの構成

```
program MAIN
use STRUCT
use PCG
use solver_ICCG_mc
implicit REAL*8 (A-H, O-Z)
real (kind=8), dimension(:), allocatable :: WK

call INPUT
call POINTER_INIT
call BOUNDARY_CELL
call CELL_METRICS
call POT_GEN

PHI= 0. d0

call solve_ICCG_mc
& ( ICELTOT, NPU, index, itemL, itemU, D,
& BFORCE, PHI, AL, AU, NCOLORTOT, PEsmpTOT,
& SMPindex, SMPindexG, EPSICCG, ITR, IER)
&
```

module STRUCT

```
module STRUCT
use omp_lib
include "precision.inc"
IC—
IC— METRICS & FLUX
integer (kind=k_int) :: ICELTOT, ICELTOTP, N
integer (kind=k_int) :: NX, NY, NZ, NXP1, NYP1, NZP1, IBNDTOT
integer (kind=k_int) :: Nxc, Nyc, Nzc
real (kind=kreal) ::
& X, Y, Z, XAREA, YAREA, ZAREA, RDx, RDY, RDZ,
& RDxZ, RDYz, RDZx, RDZy, RDZz
real (kind=kreal), dimension(:), allocatable ::
& VOLCEL, VOLIND, Vc, RW
integer (kind=k_int), dimension(:,:), allocatable ::
& XYZ, NEIBoe11
IC— BOUNDARYs
integer (kind=k_int) :: ZmaxICELTot
integer (kind=k_int), dimension(:), allocatable :: BC_INDEX, BC_MOD
integer (kind=k_int), dimension(:), allocatable :: ZmaxCEL
IC— WORK
integer (kind=k_int), dimension(:,:), allocatable :: IMX
real (kind=kreal), dimension(:,:), allocatable :: FCV
integer (kind=k_int) :: PEsmpTOT
end module STRUCT
```

ICELTOT : 要素数
ICELTOTp : = ICELTOT
N : 節点数
NX, NY, NZ : x, y, z 方向要素数
NXP1, NYP1, NZP1 : x, y, z 方向節点数
IBNDTOT : NXP1 × NYP1
XYZ (ICELTOT, 3) : 要素座標 (後述)
NEIBoe11 (ICELTOT, 6) : 隣接要素 (後述)

境界条件関連 : Z=Zmax

PEsmpTOT : スレッド数

module PCG (これまでとの相違点)

```
module PCG
integer, parameter :: N2= 256
integer :: Nmax, Nlmax, NCOLORTot, NCOLORk, NU, NL
integer :: NPL, NPU
integer :: METHOD, ORDER, METHOD
real (kind=8) :: EPSICCG
real (kind=8), dimension(:), allocatable :: D, PHI, BFORCE
real (kind=8), dimension(:), allocatable :: AL, AU
integer, dimension(:), allocatable :: IM, INU, COLORindex
integer, dimension(:), allocatable :: SMPindex, SMPindexG
integer, dimension(:), allocatable :: OLDtoNEW, NEWtoOLD
integer, dimension(:,:), allocatable :: IAL, IAU
integer, dimension(:), allocatable :: indexL, itemL
integer, dimension(:), allocatable :: indexU, itemU
end module PCG
```

NCOLORTot
COLORindex (0: NCOLORTot) 色数
各色に含まれる要素数のインデックス
(COLORindex (ico1) - COLORindex (ico1-1))

SMPindex (0: NCOLORTot+PEsmpTOT) スレッド用配列 (後述)
SMPindexG (0: PEsmpTOT)

OLDtoNEW, NEWtoOLD Coloring前後の要素番号対照表

変数表 (1/2)

配列・変数名	型	内容
D (N)	R	対角成分、(N:全メッシュ数)
BFORCE (N)	R	右辺ベクトル
PHI (N)	R	未知数ベクトル
indexL (0:N)	I	各行の非零下三角成分数 (CRS)
indexU (0:N)	I	各行の非零上三角成分数 (CRS)
NPL	I	非零下三角成分総数 (CRS)
NPU	I	非零上三角成分総数 (CRS)
itemL (NPL)	I	非零下三角成分 (列番号) (CRS)
itemU (NPU)	I	非零上三角成分 (列番号) (CRS)
AL (NPL)	R	非零下三角成分 (係数) (CRS)
AU (NPL)	R	非零上三角成分 (係数) (CRS)
NL, NU	I	各行の非零上下三角成分の最大数 (ここでは6)
INL (N)	I	各行の非零下三角成分数
INU (N)	I	各行の非零上三角成分数
IAL (NL, N)	I	各行の非零下三角成分に対応する列番号
IAU (NU, N)	I	各行の非零上三角成分に対応する列番号

変数表 (2/2)

配列・変数名	型	内容
NCOLORtot	I	入力時にはOrdering手法 (≥ 2 : MC, $=0$: CM, $=-1$: RCM, $-2 \geq$: CMRCM), 最終的には色数, レベル数が入る
COLORindex (0:NCOLORtot)	I	各色, レベルに含まれる要素数の一次元圧縮配列, COLORindex (icol-1)+1からCOLORindex (icol)までの要素がicol番目の色 (レベル) に含まれる。
NEWtoOLD (N)	I	新番号 \rightarrow 旧番号への参照配列
OLDtoNEW (N)	I	旧番号 \rightarrow 新番号への参照配列
PEsmptTOT	I	スレッド数
SMPindex (0:NCOLORtot*PEsmptTOT)	I	スレッド用補助配列 (データ依存性があるループに使用)
SMPindexG (0:PEsmptTOT)	I	スレッド用補助配列 (データ依存性が無いループに使用)

プログラムの構成

```

program MAIN
  use STRUCT
  use PCG
  use solver_ICCG_mc
  implicit REAL*8 (A-H,O-Z)
  real (Kind=8), dimension(:), allocatable :: WK
  call INPUT
  call POINTER_INIT
  call BOUNDARY_CELL
  call CELL_METRICS
  call POT_GEN
  PHI = 0.0
  &
  & call solve_ICCG_mc, NPL, NPU, indexL, itemL, indexU, itemU, D,
  & ( UCELLTOT, PHI, NL, NU, NCOLORtot, PEsmptTOT,
  & BFORCE, SMPindexG, EPSICCG, THR, TER)
  &
  &
  &

```

input

「INPUT.DAT」の読み込み

```

IC
IC*** INPUT
IC*** INPUT
IC
IC INPUT CONTROL DATA
IC
IC subroutine INPUT
  use STRUCT
  use PCG
  implicit REAL*8 (A-H,O-Z)
  character*80 CNTRL
  IC-- CNTRL, file= 'INPUT.DAT', status='unknown')
  open (11,*) file= 'INPUT.DAT', status='unknown')
  read (11,*) DZ, DZ, DZ
  read (11,*) EPSICCG
  read (11,*) PEsmptTOT
  read (11,*) NCOLORtot
  close (11)
  return
end

```

- PEsmptTOT
 - OpenMPスレッド数
- NCOLORtot
 - 色数
 - 「=0」の場合はCM
 - 「=-1」の場合はRCM
 - 「 ≤ -2 」の場合はCM-RCM

```

1.00e-02 5.00e-02 1.00e-02 DX/DY/DZ
1.00e-08 EPSICCG
16 PEsmptTOT
100 NCOLORtot

```

cell_metrics

```

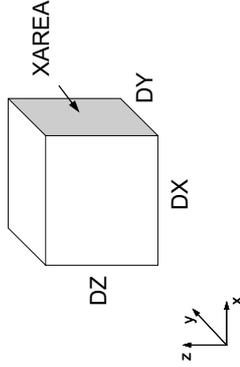
IC ***
IC *** CELL_METRICS
IC ***
IC
subroutine CELL_METRICS
  use STRUCT
  use POG
  implicit REAL*8 (A-H,O-Z)

IC-- ALLOCATE
  allocate (VOLCEL(IGELTOT))
  allocate ( RVC(IGELTOT))

IC-- VOLUME AREA PROJECTION etc.
  XAREA= DY * DZ
  YAREA= DX * DZ
  ZAREA= DX * DY
  RVZ= 1.00 / DX
  RVY= 1.00 / DY
  RVX= 1.00 / DZ
  RVZ2= 1.00 / (DX**2)
  RVY2= 1.00 / (DY**2)
  RVX2= 1.00 / (DZ**2)
  RDX= 1.00 / (0.500000X)
  RDY= 1.00 / (0.500000Y)
  RDZ= 1.00 / (0.500000Z)
  V= DX * DY * DZ
  VOLCEL= V0
  RVC = RV0
  return
end

```

計算に必要な諸パラメータ



プログラムの構成

```

program MAIN
  use STRUCT
  use POG
  use solver_ICG_mc
  implicit REAL*8 (A-H,O-Z)
  rea (kind=8), dimension(:), allocatable :: WK

  call INPUT
  call POINTER_INIT
  call BOUNDARY_CELL
  call CELL_METRICS
  call POI_GEN
  PHI= 0.00

  & call solve_ICG_mc, NPU, index, itemL, indexL, itemU, D, &
  & (IGELTOT, PH, AL, AU, NCOLORTOT, PEstTOT, &
  & BFORCE, PH, AL, AU, NCOLORTOT, PEstTOT, &
  & SMPindex, SMPindexG, EPSICG, ITR, IER)

```

poi_gen (1/9)

```

subroutine POI_GEN
  use STRUCT
  use POG
  implicit REAL*8 (A-H,O-Z)

IC-- INIT.
  m = IGETTOT
  nps = IGETTOTp

  NI= 6
  NU= 6

  alocate (BFORCE(m), D(m), PHI(m)), IAU(m),
  alocate (IU(m), INU(m), IAL(m), IAU(m), IAU(m), m))

  D = 0.00
  PHI = 0.00
  BFORCE= 0.00

  IU= 0
  INU= 0
  IAL= 0
  IAU= 0

```

配列の宣言

配列変数名	型	内容
D(N)	R	対角成分 (N:全メッシュ数)
BFORCE(N)	R	右辺ベクトル
PHI(N)	R	未知数ベクトル
indexL(O:N)	I	各行の非零下三角成分数(CRS)
indexU(O:N)	I	各行の非零上三角成分数(CRS)
NPL	I	非零下三角成分総数(CRS)
NPU	I	非零上三角成分総数(CRS)
itemL(NPL)	I	非零下三角成分列番号(CRS)
itemU(NPU)	I	非零上三角成分列番号(CRS)
AL(NPL)	R	非零下三角成分係数(CRS)
AU(NPU)	R	非零上三角成分係数(CRS)

poi_gen (2/9)

```

CONNECTIVITY
IC--
IC--
IC--
do icel= 1, IGETTOT
  icell= NEBCell(icel)
  icol2= NEBCell(icel, 1)
  icol3= NEBCell(icel, 3)
  icol4= NEBCell(icel, 4)
  icol5= NEBCell(icel, 5)
  icol6= NEBCell(icel, 6)

  if (icol6.ne.0.and.icol5.le.IGETTOT) then
    icou= IU(icel) + 1
    IAU(icou, icel)= icol5
    INU(icou, icel)= icou
  endif

  if (icol3.ne.0.and.icol4.le.IGETTOT) then
    icou= IU(icel) + icol3
    IAU(icou, icel)= icol4
    INU(icou, icel)= icou
  endif

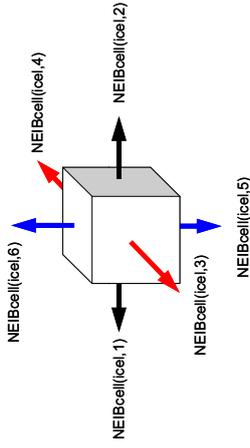
  if (icol4.ne.0.and.icol1.le.IGETTOT) then
    IAL(icou, icel)= icol1
    IAU(icou, icel)= icou
  endif

  if (icol2.ne.0.and.icol2.le.IGETTOT) then
    icou= IAU(icel) + 1
    IAU(icou, icel)= icol2
    INU(icou, icel)= icou
  endif

  if (icol4.ne.0.and.icol4.le.IGETTOT) then
    icou= IAU(icel) + 1
    IAU(icou, icel)= icol4
    INU(icou, icel)= icou
  endif

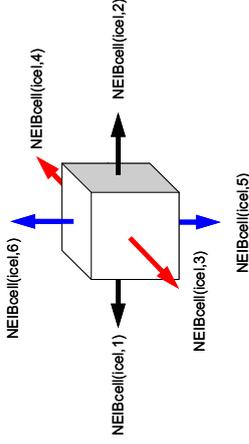
  if (icol6.ne.0.and.icol6.le.IGETTOT) then
    icou= IAU(icel) + 1
    IAU(icou, icel)= icol6
    INU(icou, icel)= icou
  endif
enddo
IC==

```



下三角成分
 NEIBcell(icol,5)= icel - NX*N
 NEIBcell(icol,3)= icel - NX
 NEIBcell(icol,1)= icel - 1

poi_gen (2/9)



上三角成分
 NEIBcell(ice1,2) = ice1 + 1
 NEIBcell(ice1,4) = ice1 + NX
 NEIBcell(ice1,6) = ice1 + NX*NY

```

CONNECTIVITY
do ice=1, NCOLORTOT
  ice1= NEIBcell(ice1,1)
  ice2= NEIBcell(ice1,2)
  ice3= NEIBcell(ice1,3)
  ice4= NEIBcell(ice1,4)
  ice5= NEIBcell(ice1,5)
  ice6= NEIBcell(ice1,6)
  if (.not. (ice1 .eq. ice2 .and. ice1 .eq. ice3 .and. ice1 .eq. ice4 .and. ice1 .eq. ice5 .and. ice1 .eq. ice6)) then
    continue
  endif
  if (.not. (ice2 .eq. ice3 .and. ice2 .eq. ice4 .and. ice2 .eq. ice5 .and. ice2 .eq. ice6)) then
    continue
  endif
  if (.not. (ice3 .eq. ice4 .and. ice3 .eq. ice5 .and. ice3 .eq. ice6)) then
    continue
  endif
  if (.not. (ice4 .eq. ice5 .and. ice4 .eq. ice6)) then
    continue
  endif
  if (.not. (ice5 .eq. ice6)) then
    continue
  endif
  !locate (OLDINDEX(ICE1TOT), NEWCOLID(ICE1TOT))
  !locate (COLINDEX(ICE1TOT), NEWCOLID(ICE1TOT))
  continue
  write (*, '(a, i8.a)') 'You have', ICE1TOT, ' elements.'
  write (*, '(a)') 'How many colors do you need?'
  write (*, '(a)') 'The answer is', NCOLORTOT, ' because there are', ICE1TOT, ' colors.'
  write (*, '(a, i8)') '#COLOR =', NCOLORTOT, ' CM if #COLOR .eq. 0'
  write (*, '(a)') 'CM if #COLOR .eq. 1'
  write (*, '(a)') 'RCM if #COLOR .eq. 1'
  write (*, '(a)') 'CMRDM if #COLOR .le. 2'
  write (*, '(a)') ' '
  if (NCOLORTOT .eq. 0) then
    call IC (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  if (NCOLORTOT .eq. 0) then
    call CM (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  if (NCOLORTOT .eq. 1) then
    call RCM (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  if (NCOLORTOT .lt. 2) then
    call CMRDM (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  write (*, '(a, i8, //)') '## FINAL COLOR NUMBER', NCOLORTOT
enddo

```

poi_gen (3/9)

並べ替えの実施：
 NCOLORTOT > 1 : Multicolor
 NCOLORTOT = 0 : CM
 NCOLORTOT = 1 : RCM
 NCOLORTOT < 1 : CM-RCM

```

!C
!C | MULTI-COLORING |
!C
!C==
do ice=1, NCOLORTOT
  !locate (OLDINDEX(ICE1TOT), NEWCOLID(ICE1TOT))
  !locate (COLINDEX(ICE1TOT), NEWCOLID(ICE1TOT))
  continue
  write (*, '(a, i8.a)') 'You have', ICE1TOT, ' elements.'
  write (*, '(a)') 'How many colors do you need?'
  write (*, '(a)') 'The answer is', NCOLORTOT, ' because there are', ICE1TOT, ' colors.'
  write (*, '(a, i8)') '#COLOR =', NCOLORTOT, ' CM if #COLOR .eq. 0'
  write (*, '(a)') 'CM if #COLOR .eq. 1'
  write (*, '(a)') 'RCM if #COLOR .eq. 1'
  write (*, '(a)') 'CMRDM if #COLOR .le. 2'
  write (*, '(a)') ' '
  if (NCOLORTOT .eq. 0) then
    call IC (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  if (NCOLORTOT .eq. 0) then
    call CM (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  if (NCOLORTOT .eq. 1) then
    call RCM (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  if (NCOLORTOT .lt. 2) then
    call CMRDM (ICE1TOT, NL, INL, IAL, INL, IAU, INCOLORTOT, COLORINDEX, NEIBCOLID, OLDINDEX) &
  endif
  write (*, '(a, i8, //)') '## FINAL COLOR NUMBER', NCOLORTOT
enddo

```

poi_gen (4/9)

各色内の要素数：
 COLORINDEX(ice) - COLORINDEX(ice-1)
 同じ色内の要素は依存性が無いため、
 並列に計算可能 ⇒ OpenMP 適用
 これを更に「PEsmptOT」で割って
 「SMPindex」に割り当てる。

前処理で使用

```

!locate (SMPindex(0:PEsmptOT+NCOLORTOT))
do SMPindex=0, NCOLORTOT
  nr= COLORINDEX(ice) - COLORINDEX(ice-1)
  nr= nr / PEsmptOT
  nr= nr - PEsmptOT*nr
  do ip=1, PEsmptOT
    if (.not. (ip .le. nr)) then
      SMPindex((ice-1)*PEsmptOT+ip)= nr + 1
    else
      SMPindex((ice-1)*PEsmptOT+ip)= nr
    endif
  enddo
enddo

do ice=1, NCOLORTOT
  do ip=1, PEsmptOT
    j1= (ice-1)*PEsmptOT + ip
    j2= j1 - 1
    SMPindex(j1)= SMPindex(j0) + SMPindex(j1)
  enddo
enddo

!locate (SMPindex(0:PEsmptOT))
SMPindexG=0
nr= ICE1TOT / PEsmptOT
nr= ICE1TOT - nr*PEsmptOT
do ip=1, PEsmptOT
  SMPindexG(ip)= nr
  if (.not. (ip .le. nr)) SMPindexG(ip)= nr + 1
enddo

do ip=1, PEsmptOT
  SMPindexG(ip)= SMPindexG(ip-1) + SMPindexG(ip)
enddo

```

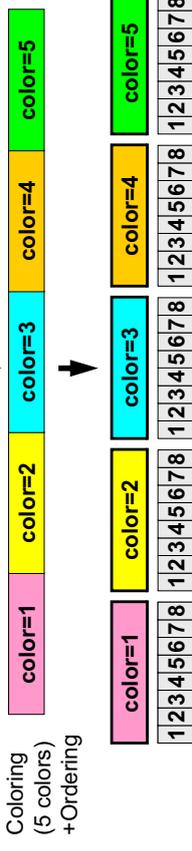
SMPindex: 前処理向け

```

do ice=1, NCOLORTOT
  !omp parallel do ...
  do ip=1, PEsmptOT
    ip1= (ice-1)*PEsmptOT+ip
    do i= SMPindex(ip1-1)+1, SMPindex(ip1)
      (... )
    enddo
  enddo
!omp end parallel do
enddo

```

Initial Vector



- 5色, 8スレッドの例
- 同じ「色」に属する要素は独立⇒並列計算可能
- 色の順番に並び替え

不完全修正コレスキー分解：並列版

$$d_i = \left(a_{ii} - \sum_{k=1}^{i-1} a_{ik}^2 \cdot d_k \right)^{-1} = l_{ii}^{-1}$$

$W(i, DD):$ d_i
 $D(i):$ a_{ii}
 $IAL(j, i):$ k
 $AL(j, i):$ a_{ik}

```

do i= 1, NCOLORTot
!$omp parallel do private(ip, ip1, i, VAL, k)
do ip= i, PEsmpTOT
ip1= (i-1)*PEsmpTOT + ip
do i= SMPindex(ip1-1)+1, SMPindex(ip1)
VAL= D(i)
do k= indexL(i-1)+1, indexL(i)
VAL= VAL - (AL(k)**2) * W(itemL(k), DD)
enddo
W(i, DD)= 1. d0/VAL
enddo
enddo
!$omp end parallel do

```

privateに注意。

行列ベクトル積

依存性が無い⇒独立に計算可能⇒SMPindexG使用

```

!$omp parallel do private(ip, i, VAL, k)
do ip= 1, PEsmpTOT
do i = SMPindexG(ip-1)+1, SMPindexG(ip)
VAL= D(i) * X(i)
do k= indexL(i-1)+1, indexL(i)
VAL= VAL + AL(k) * X(itemL(k))
enddo
do k= indexU(i-1)+1, indexU(i)
VAL= VAL + AU(k) * X(itemU(k))
enddo
W(i, R)= B(i) - VAL
enddo
enddo
!$omp end parallel do

```

solve_ICCG_mc(3/6)

```

iC |r0|= |b| - |A| |xini| |
iC |
!$omp parallel do private(ip, i, VAL, k)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
VAL= D(i) * X(itemL(i-1))+1, indexL(i)
do k= indexL(i-1)+1, indexL(i)
VAL= VAL + AL(k) * X(itemL(k))
enddo
do k= indexU(i-1)+1, indexU(i)
VAL= VAL + AU(k) * X(itemU(k))
enddo
W(i, R)= B(i) - VAL
enddo
enddo
!$omp end parallel do
BNRM2= 0.000
!$omp parallel do private(ip, i, VAL, k) reduction(←:BNRM2)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
BNRM2 = BNRM2 + B(i) **2
enddo
enddo
!$omp end parallel do

```

```

Compute r(0)= b-[A]x(0)
for i= 1, 2, ...
solve [M]z^(i-1)= r^(i-1)
p_{i-1}= r^(i-1) z^(i-1)
if i=1
p^(1)= z(0)
else
beta_{i-1}= p_{i-1}/p_{i-2}
p^(i)= z^(i-1) + beta_{i-1} p^(i-1)
endif
q^(i)= [A]p^(i)
alpha_i = p_{i-1}/p^(i) q^(i)
x^(i)= x^(i-1) + alpha_i p^(i)
r^(i)= r^(i-1) - alpha_i q^(i)
check convergence |r|
end

```

solve_ICCG_mc(3/6)

```

iC |r0|= |b| - |A| |xini| |
iC |
!$omp parallel do private(ip, i, VAL, k)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
VAL= D(i) * X(itemL(i-1))+1, indexL(i)
do k= indexL(i-1)+1, indexL(i)
VAL= VAL + AL(k) * X(itemL(k))
enddo
do k= indexU(i-1)+1, indexU(i)
VAL= VAL + AU(k) * X(itemU(k))
enddo
W(i, R)= B(i) - VAL
enddo
enddo
!$omp end parallel do
BNRM2= 0.000
!$omp parallel do private(ip, i, VAL, k) reduction(←:BNRM2)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
BNRM2 = BNRM2 + B(i) **2
enddo
enddo
!$omp end parallel do

```

```

Compute r(0)= b-[A]x(0)
for i= 1, 2, ...
solve [M]z^(i-1)= r^(i-1)
p_{i-1}= r^(i-1) z^(i-1)
if i=1
p^(1)= z(0)
else
beta_{i-1}= p_{i-1}/p_{i-2}
p^(i)= z^(i-1) + beta_{i-1} p^(i-1)
endif
q^(i)= [A]p^(i)
alpha_i = p_{i-1}/p^(i) q^(i)
x^(i)= x^(i-1) + alpha_i p^(i)
r^(i)= r^(i-1) - alpha_i q^(i)
check convergence |r|
end

```

内積: SMPindexG使用, reduction

```

BNRM2= 0.0D0
!$omp parallel do private(ip,i) reduction(+:BNRM2)
do ip= 1, PEsmpTOT
do i = SMPindexG(ip-1)+1, SMPindexG(ip)
BNRM2 = BNRM2 + B(i) **2
enddo
enddo
!$omp end parallel do

```

solve_ICCG_mc(4/6)

```

ITR= N
do L= 1, ITR
do i= 1, 2, ...
!$omp parallel do private(ip,i)
do ip= (i-1)*PEsmpTOT + 1, SMPindexG(ip)
W(i,Z)= W(i,R)
enddo
!$omp end parallel do
do i= 1, NCOLTOT
!$omp parallel do private(ip,i,WAL,j)
do ip= 1, PEsmpTOT
do j= (i-1)*PEsmpTOT + 1, SMPindex(ip)
WAL= W(i,Z)
do k= 1, INU(i)
WAL= WAL - AL(k) * W(iAU(j),Z)
enddo
W(i,Z)= WAL * W(i,DD)
enddo
enddo
!$omp end parallel do
do i= NCOLTOT, 1, -1
!$omp parallel do private(ip,i,SW,j)
do ip= 1, PEsmpTOT
do j= (i-1)*PEsmpTOT + 1, SMPindex(ip)
SW= SMPindex(ip-1)+1, INU(i)
do k= 1, INU(i)
SW= SW + AU(j,k) * W(iAU(j),Z)
enddo
W(i,Z)= W(i,Z) - W(i,DD) * SW
enddo
enddo
!$omp end parallel do
enddo

```

```

Compute  $x^{(0)} = b - [A]x^{(0)}$ 
for  $i = 1, 2, \dots$ 
  solve  $[M]z^{(i-1)} = r^{(i-1)}$ 
   $\rho_{i-1} = x^{(i-1)} z^{(i-1)}$ 
  if  $i=1$ 
     $p^{(1)} = z^{(0)}$ 
  else
     $\rho_{i-1} = \rho_{i-1} / \rho_{i-2}$ 
     $p^{(i)} = z^{(i-1)} + \beta_{i-1} p^{(i-1)}$ 
  endif
   $q^{(i)} = [A]p^{(i)}$ 
   $\alpha_i = \rho_{i-1} / p^{(i)} q^{(i)}$ 
   $x^{(i)} = x^{(i-1)} + \alpha_i p^{(i)}$ 
   $r^{(i)} = r^{(i-1)} - \alpha_i q^{(i)}$ 
  check convergence | r |
end

```

solve_ICCG_mc(4/6)

```

ITR= N
do L= 1, ITR
do i= 1, 2, ...
!$omp parallel do private(ip,i)
do ip= (i-1)*PEsmpTOT + 1, SMPindexG(ip)
W(i,Z)= W(i,R)
enddo
!$omp end parallel do
do i= 1, NCOLTOT
!$omp parallel do private(ip,i,WAL,k)
do ip= 1, PEsmpTOT
do k= (i-1)*PEsmpTOT + 1, SMPindex(ip)
WAL= W(i,Z)
do l= indexL(i-1)+1, indexU(i)
WAL= WAL - AL(k) * W(iTemL(k),Z)
enddo
W(i,Z)= WAL * W(i,DD)
enddo
enddo
!$omp end parallel do
do i= NCOLTOT, 1, -1
!$omp parallel do private(ip,i,SW,k)
do ip= 1, PEsmpTOT
do k= (i-1)*PEsmpTOT + 1, SMPindex(ip)
SW= SMPindex(ip-1)+1, indexU(i)
do l= indexL(i-1)+1, indexU(i)
SW= SW + AU(k) * W(iTemU(k),Z)
enddo
W(i,Z)= W(i,Z) - W(i,DD) * SW
enddo
enddo
!$omp end parallel do
enddo

```

```

Compute  $x^{(0)} = b - [A]x^{(0)}$ 
for  $i = 1, 2, \dots$ 
  solve  $[M]z^{(i-1)} = r^{(i-1)}$ 
   $\rho_{i-1} = x^{(i-1)} z^{(i-1)}$ 
  if  $i=1$ 
     $p^{(1)} = z^{(0)}$ 
  else
     $\rho_{i-1} = \rho_{i-1} / \rho_{i-2}$ 
     $p^{(i)} = z^{(i-1)} + \beta_{i-1} p^{(i-1)}$ 
  endif
   $q^{(i)} = [A]p^{(i)}$ 
   $\alpha_i = \rho_{i-1} / p^{(i)} q^{(i)}$ 
   $x^{(i)} = x^{(i-1)} + \alpha_i p^{(i)}$ 
   $r^{(i)} = r^{(i-1)} - \alpha_i q^{(i)}$ 
  check convergence | r |
end

```

solve_ICCG_mc(4/6)

```

ITR= N
do L= 1, ITR
do i= 1, 2, ...
!$omp parallel do private(ip,i)
do ip= (i-1)*PEsmpTOT + 1, SMPindexG(ip)
W(i,Z)= W(i,R)
enddo
!$omp end parallel do
do i= 1, NCOLTOT
!$omp parallel do private(ip,i,WAL,k)
do ip= 1, PEsmpTOT
do k= (i-1)*PEsmpTOT + 1, SMPindex(ip)
WAL= W(i,Z)
do l= indexL(i-1)+1, indexU(i)
WAL= WAL - AL(k) * W(iTemL(k),Z)
enddo
W(i,Z)= WAL * W(i,DD)
enddo
enddo
!$omp end parallel do
do i= NCOLTOT, 1, -1
!$omp parallel do private(ip,i,SW,k)
do ip= 1, PEsmpTOT
do k= (i-1)*PEsmpTOT + 1, SMPindex(ip)
SW= SMPindex(ip-1)+1, indexU(i)
do l= indexL(i-1)+1, indexU(i)
SW= SW + AU(k) * W(iTemU(k),Z)
enddo
W(i,Z)= W(i,Z) - W(i,DD) * SW
enddo
enddo
!$omp end parallel do
enddo

```

$$(M)\{z\} = (LDL^T)\{z\} = \{r\}$$

前進代入
Forward Substitution

$$(L)\{z\} = \{r\}$$

後退代入
Backward Substitution

$$(D^T)\{z\} = \{z\}$$

前進代入 : SMPindex使用

```

do ic= 1, NCOLORTot
!$omp parallel do private(ip, ip1, i, WVAL, k)
do ip= 1, PEsmpTOT
ip1= (ic-1)*PEsmpTOT + ip
do i= SMPindex(ip1-1)+1, SMPindex(ip1)
WVAL= W(i, Z)
do k= indexL(i-1)+1, indexL(i)
WVAL= WVAL - AL(k) * W(indexL(k), Z)
enddo
W(i, Z)= WVAL * W(i, DD)
enddo
enddo
!$omp end parallel do
enddo

```

solve_ICCG_mc(5/6)

```

ic | (p) = (z) if ITER=1
ic | BETA= RHO / RHO1 otherwise
ic ==
if (L.eq.1) then
!$omp parallel do private(ip, i)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
W(i, P)= W(i, Z)
enddo
enddo
else
!$omp end parallel do
!$omp parallel do private(ip, i, VAL, k)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
VAL= D(i)*W(i, P)
do k= indexL(i-1)+1, indexL(i)
VAL= VAL + AL(k)*W(indexL(k), P)
enddo
do k= indexU(i-1)+1, indexU(i)
VAL= VAL + AU(k)*W(indexU(k), P)
enddo
W(i, D)= VAL
enddo
enddo
!$omp end parallel do
enddo

```

```

Compute  $r^{(0)} = b - [A]x^{(0)}$ 
for i= 1, 2, ...
solve  $[M]z^{(i-1)} = r^{(i-1)}$ 
 $\rho_{i-1} = r^{(i-1)} z^{(i-1)}$ 
if i=1
 $p^{(1)} = z^{(0)}$ 
else
 $\beta_{i-1} = \rho_{i-1} / \rho_{i-2}$ 
 $p^{(i)} = z^{(i-1)} + \beta_{i-1} p^{(i-1)}$ 
endif
 $q^{(i)} = [A]p^{(i)}$ 
 $\alpha_i = \rho_{i-1} / p^{(i)} q^{(i)}$ 
 $x^{(i)} = x^{(i-1)} + \alpha_i p^{(i)}$ 
 $r^{(i)} = r^{(i-1)} - \alpha_i q^{(i)}$ 
check convergence | r |
end

```

solve_ICCG_mc(5/6)

```

ic | (p) = (z) if ITER=1
ic | BETA= RHO / RHO1 otherwise
ic ==
if (L.eq.1) then
!$omp parallel do private(ip, i)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
W(i, P)= W(i, Z)
enddo
enddo
else
!$omp end parallel do
!$omp parallel do private(ip, i)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
W(i, P)= W(i, Z) + BETA*W(i, P)
enddo
enddo
!$omp end parallel do
enddo

```

```

Compute  $r^{(0)} = b - [A]x^{(0)}$ 
for i= 1, 2, ...
solve  $[M]z^{(i-1)} = r^{(i-1)}$ 
 $\rho_{i-1} = r^{(i-1)} z^{(i-1)}$ 
if i=1
 $p^{(1)} = z^{(0)}$ 
else
 $\beta_{i-1} = \rho_{i-1} / \rho_{i-2}$ 
 $p^{(i)} = z^{(i-1)} + \beta_{i-1} p^{(i-1)}$ 
endif
 $q^{(i)} = [A]p^{(i)}$ 
 $\alpha_i = \rho_{i-1} / p^{(i)} q^{(i)}$ 
 $x^{(i)} = x^{(i-1)} + \alpha_i p^{(i)}$ 
 $r^{(i)} = r^{(i-1)} - \alpha_i q^{(i)}$ 
check convergence | r |
end

```

solve_ICCG_mc(6/6)

```

ic | ALPHA= RHO / (p) | (q)
ic ==
G1= 0.40
!$omp parallel do private(ip, i) reduction(+:G1)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
G1= G1 + W(i, P)*W(i, Q)
enddo
enddo
ALPHA= RHO / G1
ic ==
ic | (x) = (x) + ALPHA*(p)
ic | (r) = (r) - ALPHA*(q)
ic ==
!$omp parallel do private(ip, i)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
W(i, P)= W(i, P) + ALPHA * W(i, Q)
enddo
enddo
DNRHO= 0.40
!$omp parallel do private(ip, i) reduction(+:DNRHO)
do ip= 1, PEsmpTOT
do i= SMPindex(ip-1)+1, SMPindex(ip)
DNRHO= DNRHO + W(i, P)*W(i, Q)
enddo
enddo
!$omp end parallel do
enddo

```

```

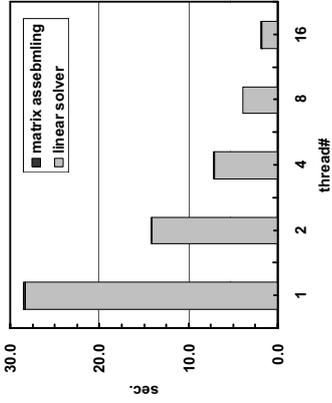
Compute  $r^{(0)} = b - [A]x^{(0)}$ 
for i= 1, 2, ...
solve  $[M]z^{(i-1)} = r^{(i-1)}$ 
 $\rho_{i-1} = r^{(i-1)} z^{(i-1)}$ 
if i=1
 $p^{(1)} = z^{(0)}$ 
else
 $\beta_{i-1} = \rho_{i-1} / \rho_{i-2}$ 
 $p^{(i)} = z^{(i-1)} + \beta_{i-1} p^{(i-1)}$ 
endif
 $q^{(i)} = [A]p^{(i)}$ 
 $\alpha_i = \rho_{i-1} / p^{(i)} q^{(i)}$ 
 $x^{(i)} = x^{(i-1)} + \alpha_i p^{(i)}$ 
 $r^{(i)} = r^{(i-1)} - \alpha_i q^{(i)}$ 
check convergence | r |
end

```


計算時間 (MC=2)

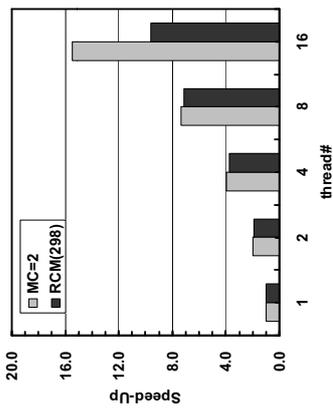
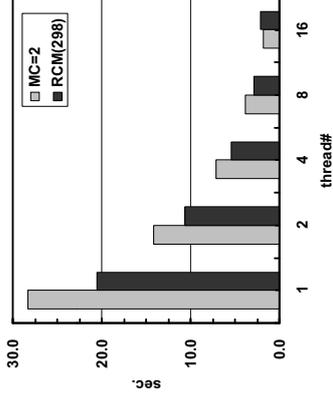
以降, 「linear solver」の計算時間のみ問題とする

matrix assembling: poi-genの後半: 並列化
linear solver: CG



スケールABILITY

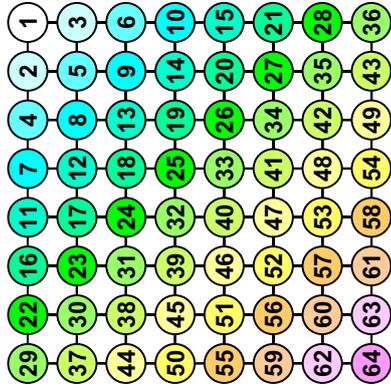
反復回数: MC (2色) : 333回, RCM (298レベル) : 224回



- 反復回数ではRCMの方が少ないのに, コア数が増えると, 計算時間が逆転する
- MC=2は良好なスケールABILITY (16コアで15.5)

原因

- 色数が多い (298対2) ⇒ 同期オーバーヘッド
- スレッドごとの不可不均衡
 - RCMでは要素数の少ない「レベル」が必ず存在する



原因はこのようなループにある

色数増加 ⇒ 同期オーバーヘッド増加
必ずある「色」の計算が終わってから次の「色」に行く

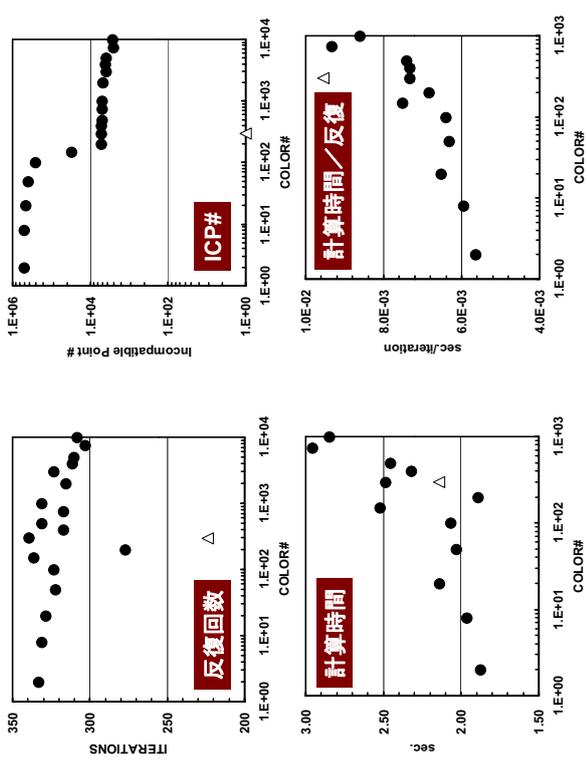
```

do ic= 1, NCOLORTot
!$omp parallel do private(ip, ip1, i, WVAL, k)
do ip= 1, PEsmpTOT
  ip1= (ic-1)*PEsmpTOT + ip
  do i= SMPindex(ip1-1)+1, SMPindex(ip1)
    WVAL= W(i, Z)
    do k= indexL(i-1)+1, indexL(i)
      WVAL= WVAL - AL(k) * W(itemL(k), Z)
    enddo
    W(i, Z)= WVAL * W(i, DD)
  enddo
enddo
!$omp end parallel do
enddo
  
```

MCとRCMの比較

- MC
 - 並列性高い, 負荷分散も良い(そのように設定されている)
 - 特に色数少ないと反復回数多い
 - 色数を増やす, 反復回数減るが同期オーバーヘッドの影響
- RCM
 - 収束は早い, レベル数が多く, 同期オーバーヘッドの影響受けやすく, コア数が増えると不利
 - 負荷分散もいま一つ
- **反復回数少なくて同期オーバーヘッドの影響が少ない方法が無いものか?**
 - **色数が少なくて, かつ反復回数が少ないという都合の良い方法**

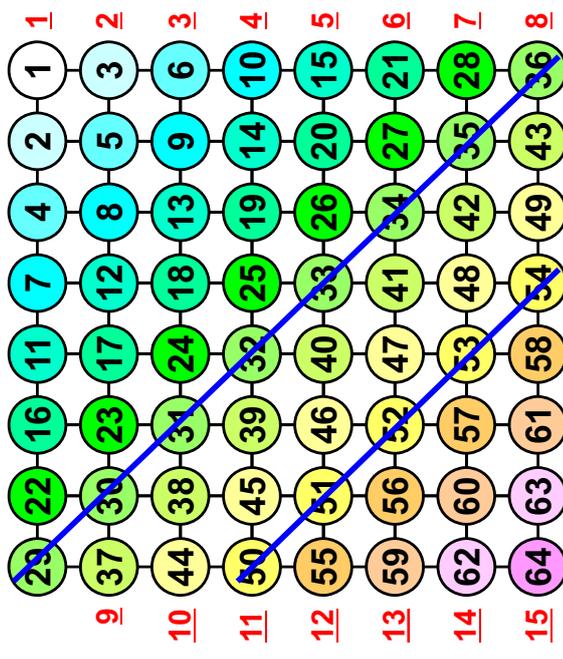
16コアにおける結果(●:MC, △:RCM)



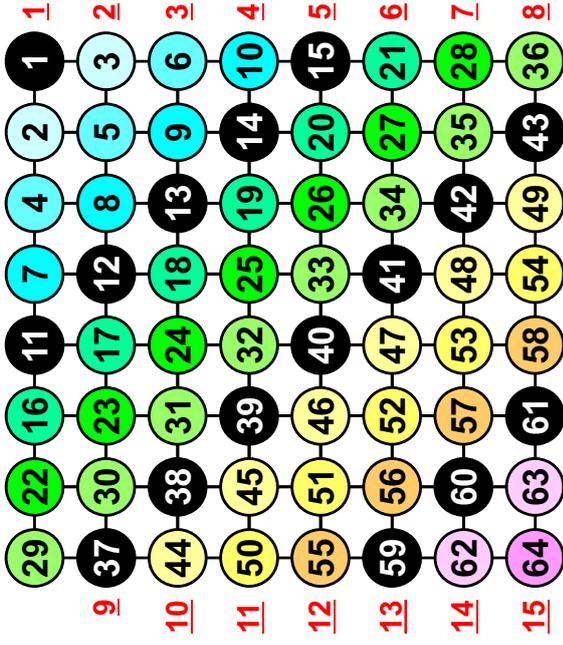
解決策: CM-RCM

- RCM + Cyclic Multicoloring [土井, 襲田, 鷺尾他]
- 手順
 - まずRCMを施す
 - Cyclic Multicoloring (CM) の色数を決める (N_c)
 - RCMの1番目, (N_c+1) 番目, $(2N_c+1)$ 番目...のレベルに属する要素を「1」色に分類する
 - RCMのk番目, (N_c+k) 番目, $(2N_c+k)$ 番目...のレベルに属する要素を「k」色に分類する
 - 「k」が「 N_c 」に達して, 要素が「1~ N_c 」で色付けされたら完了
 - あとはMCのときと同じように, 色の順番に再番号付
 - RCMの各レベルに対して「 N_c 」のサイクルで再色付けを実施している
 - **もし同じ色の要素の中に依存性が見つかったら, $N_c=N_c+1$ として最初からやり直し(ここは少し原始的)**

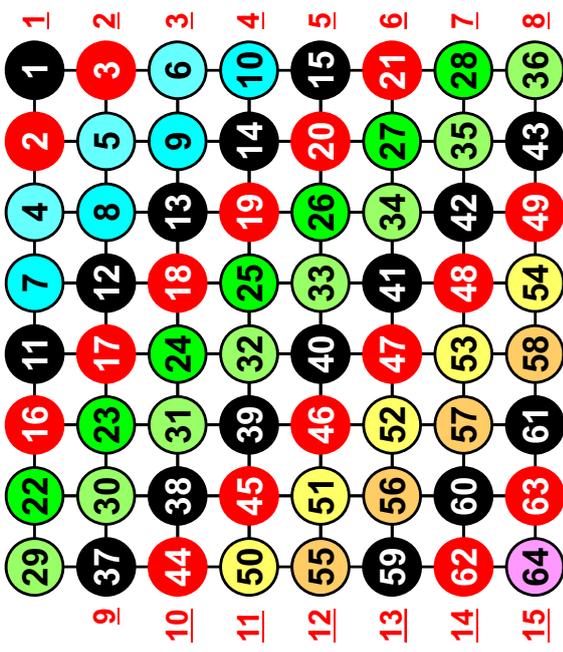
RCM



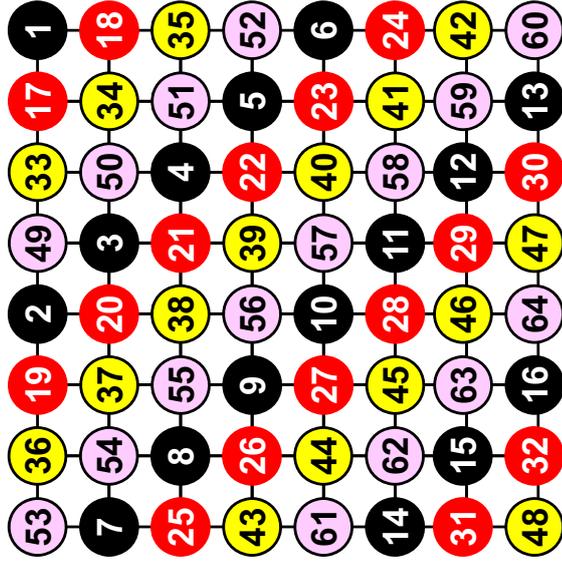
Nc=4, k=1: 1,5,9,13レベルを選択



Nc=4, k=2: 2,6,10,14レベルを選択

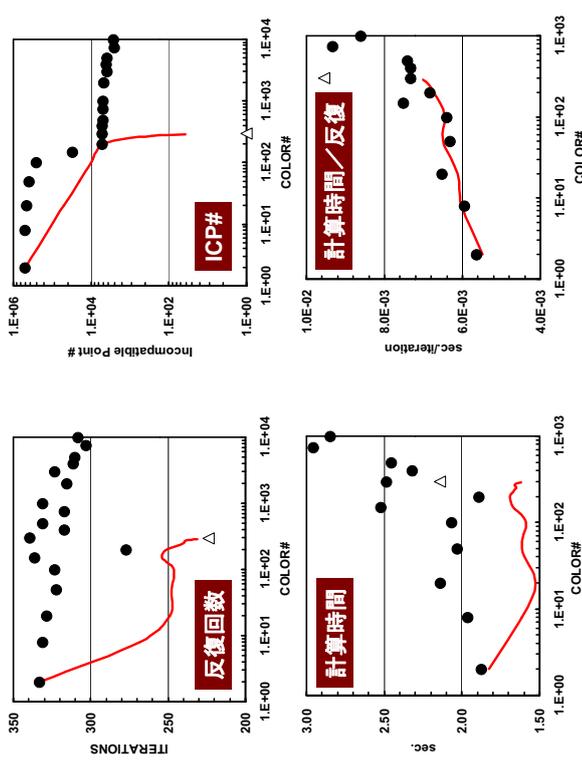


CM-RCM(Nc=4): 「色」の順番に再並替



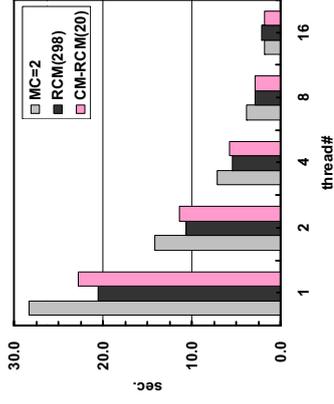
k=1: 16
k=2: 16
k=3: 16
k=4: 16

16コアにおける結果 (●: MC, △: RCM)



スケーラビリティ

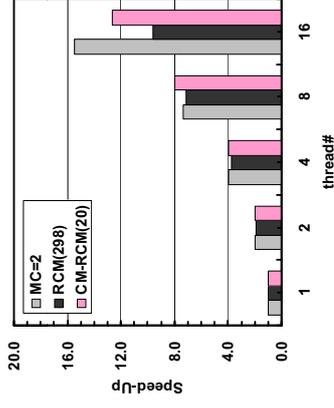
反復回数: MC(2色): 333回, RCM(298レベル): 224回
 CM-RCM(Nc=20): 249回



16 threads

MC(2): 1.83 sec.

CM-RCM(20): 1.79 sec.



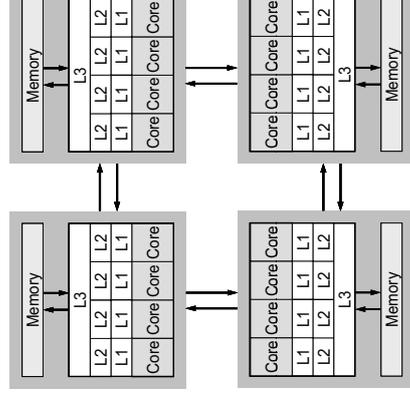
CM-RCM

- 少ない色数(Nc)で良好な収束を得られる
- 計算効率も良い
- 実行方法
 - INPUT.DATで「NCOLORtot=-Nc」とする
 - L2, L3で有効なオプション(既にL2でも使っていた)
- 実装は「cmrcm.f」を参照ください(本日は説明は省略)

- L2-solへのOpenMPの実装
- Hitachi SR11000/J2での実行
 - 計算結果
 - CM-RCMオーダリング
- **T2K(東大)での実行**
- T2K(東大)での性能向上への道
 - NUMA Control
 - First Touch
 - データ再配置

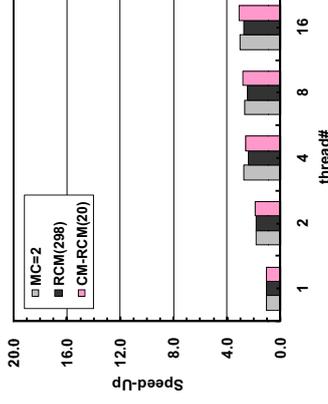
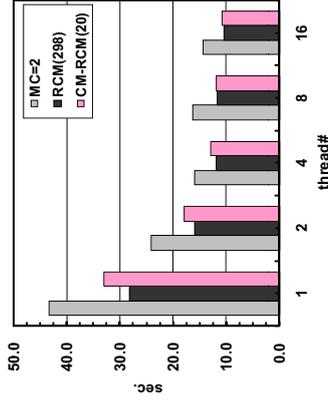
計算結果

- T2K(東大) 1ノード(4ソケット, 16コア)
- 100³要素



スケーラビリティ: CASE-0

反復回数: MC(2色): 333回, RCM(298レベル): 224回
CM-RCM(Nc=20): 249回



実行用データ

INPUT.DAT

```
1.00e-00 1.00e-00 1.00e-00
1.0e-08
16
100
DX/DX/DZ
EPSICCG
PEmpTOT (固定)
NCOLORtot
```

x0.sh

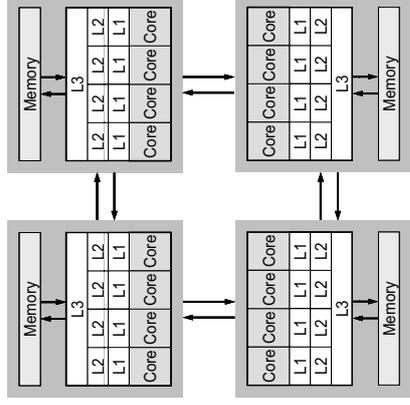
```
#@$-r test
#@$-g tutorial
#@$-N 1
#@$-J T1
#@$-e err
#@$-o test.lst
#@$-lM 28GB
#@$-lE 00:10:00
#@$-s /bin/sh
#@$
cd $PBS_O_WORKDIR
export OMP_NUM_THREADS= 16
mpirun ./I3-sol
exit
```

性能向上への道

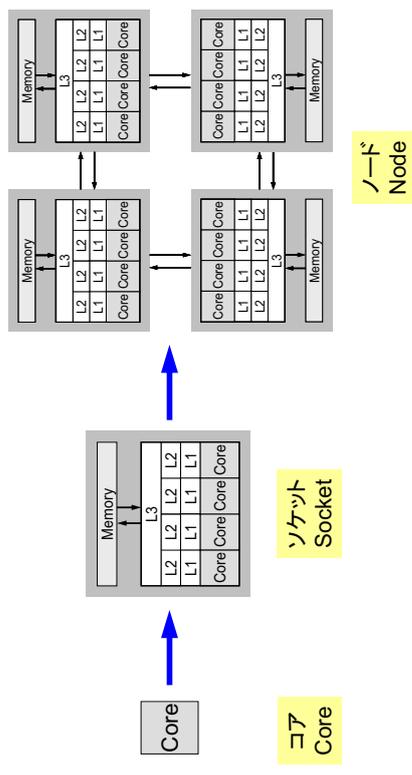
- L2-solへのOpenMPの実装
- Hitachi SR11000/J2での実行
 - 計算結果
 - CM-RCMオーダリング
- T2K(東大)での実行
- **T2K(東大)での性能向上への道**
 - **NUMA Control**
 - First Touch
 - データ再配置
- CASE-0 初期状態
- CASE-1a NUMA control (policy-a)
- CASE-1b NUMA control (policy-b)
- CASE-2a First-touch (a)
- CASE-2b First-touch (b)
- CASE-3a First-touch+データ再配置(a)
- CASE-3b First-touch+データ再配置(b)
- CASE-4a データ再配置(a)
- CASE-4b データ再配置(b)

Quad Core Opteron: NUMA Architecture

- AMD Quad Core Opteron 2.3GHz
 - Quad Coreのソケット × 4 ⇒ 1ノード(16コア)
- 各ソケットがローカルにメモリを持っている
 - cc-NUMA: cache-coherent-Non-Uniform Memory Access
 - できるだけローカルのメモリをアクセスして計算するようなプログラムミング, データ配置, 実行時制御 (numactl) が必要



NUMA Architecture



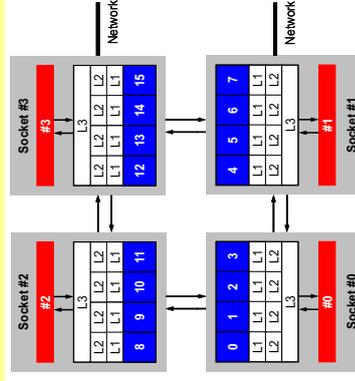
numactl

- NUMA(Non Uniform Memory Access) 向けのメモリ割付のためのコマンドライン: Linuxでサポート
- T2K(東大)でも実績有り[Nakajima, K. 2008 (IEEE Cluster 2008)]

```
>$ numactl --show
policy: default
preferred node: current
physcpubind: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
cpubind: 0 1 2 3
nodebind: 0 1 2 3
membind: 0 1 2 3
```

numactl --show

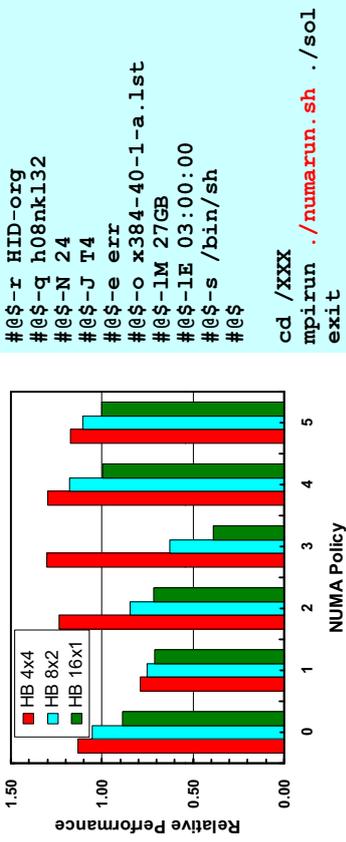
```
>$ numactl --show
policy: default
preferred node: current
physcpubind: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
cpubind: 0 1 2 3
nodebind: 0 1 2 3
membind: 0 1 2 3
```



例：numactlの影響

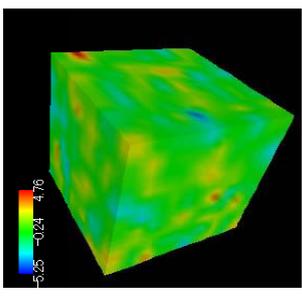
- T2K(東大), 有限要素法アプリケーション
- POLICY=0:何も指定しない場合
- 相対性能:大きいほど良い
 - Flat MPIに対する相対性能

状況によって、最適な組み合わせは実は異なる(後述)

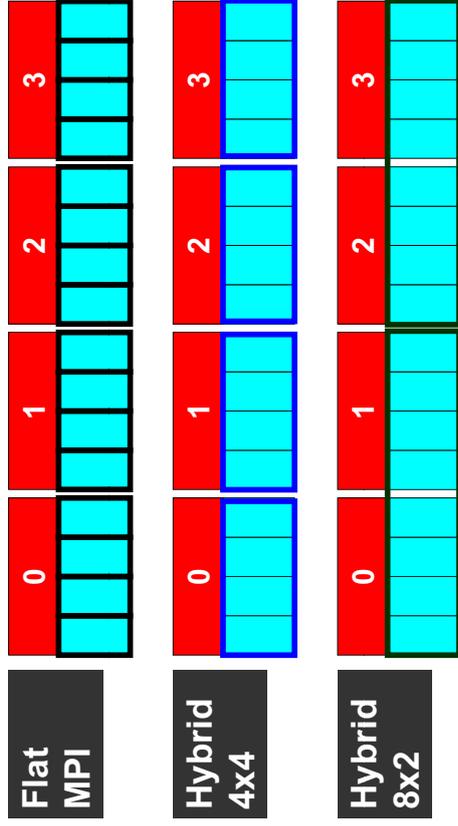


Target Application

- 3D Elastic Problems with Heterogeneous Material Property
 - $E_{max}=10^3, E_{min}=10^{-3}, \nu=0.25$
 - generated by "sequential Gauss" algorithm for geo-statistics [Deutsch & Journel, 1998]
 - 128³ tri-linear hexahedral elements, 6,291,456 DOF
- (SGS+GPBiCG) Iterative Solvers
 - Symmetric Gauss-Seidel
 - Original Block Jacobi, HID
- T2K/Tokyo
 - 512 cores (32 nodes)
- FORTARN90 (Hitachi) + MPI
 - Flat MPI, Hybrid (4x4, 8x2)
- Effect of NUMA Control



Flat MPI, Hybrid (4x4, 8x2)



numarun.shの中身

NUMA Policy	HB 4x4	HB 8x2	HB 16x1
0	~0.8	~1.1	~1.0
1	~0.9	~1.0	~0.9
2	~1.0	~1.0	~0.9
3	~1.0	~1.0	~0.9
4	~1.0	~1.0	~0.9
5	~1.0	~1.0	~0.9

```

Policy:1
#!/bin/bash
MYRANK=$MXMPI_ID MPIプロセス番号
MYVAL=$(expr $MYVAL / 4)
SOCKET=$(expr $MYVAL % 4)
numactl --cpunodebind=$SOCKET --interleave=all $@

Policy:2
#!/bin/bash
MYRANK=$MXMPI_ID
MYVAL=$(expr $MYRANK / 4)
SOCKET=$(expr $MYVAL % 4)
numactl --cpunodebind=$SOCKET --interleave=$SOCKET $@

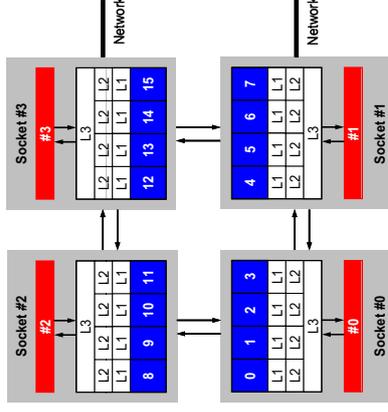
Policy:3
#!/bin/bash
MYRANK=$MXMPI_ID
MYVAL=$(expr $MYRANK / 4)
SOCKET=$(expr $MYVAL % 4)
numactl --cpunodebind=$SOCKET --interleave=$SOCKET $@

Policy:4
#!/bin/bash
MYRANK=$MXMPI_ID
MYVAL=$(expr $MYRANK / 4)
SOCKET=$(expr $MYVAL % 4)
numactl --cpunodebind=$SOCKET --membind=$SOCKET $@

Policy:5
#!/bin/bash
MYRANK=$MXMPI_ID
MYVAL=$(expr $MYVAL % 4)
SOCKET=$(expr $MYVAL % 4)
numactl --localalloc $@
  
```

本ケースにおける:numactlのポリシー

- policy-a
 - --localalloc
 - ローカルなソケットのメモリ使用
- policy-b
 - --interleave= all
 - ノード上のメモリをinterleave



本ケースにおける:numactlのポリシー

x5.sh Policy-a

```
#@$-r test
#@$-q tutorial
#@$-N 1
#@$-J T1
#@$-e err
#@$-o test.lst
#@$-IM 28GB
#@$-IE 00:10:00
#@$-s /bin/sh
#@$

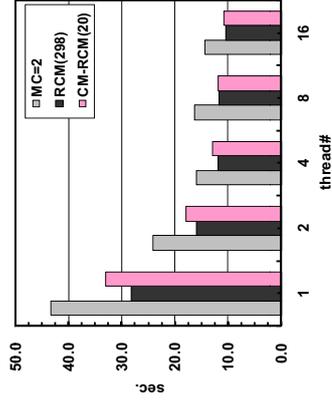
cd $PBS_O_WORKDIR
export OMP_NUM_THREADS= 16
mpirun numactl --localalloc ./sol
exit
```

x6.sh Policy-b

```
cd $PBS_O_WORKDIR
export OMP_NUM_THREADS= 16
mpirun numactl --interleave=all ./sol
exit
```

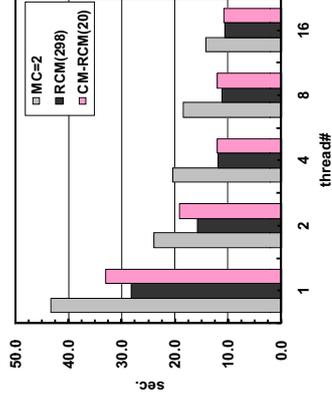
スケラビリティ:numactlの効果

反復回数: MC (2色) : 333回, RCM(298レベル) : 224回
 CM-RCM(Nc=20) : 249回
CASE-1a: あまり変わらない



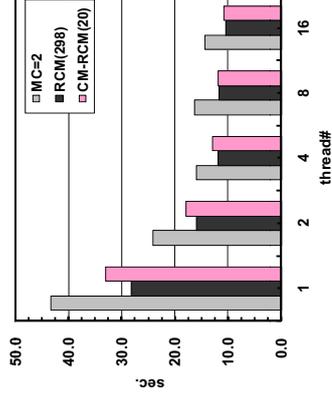
CASE-0

CASE-1a --localalloc



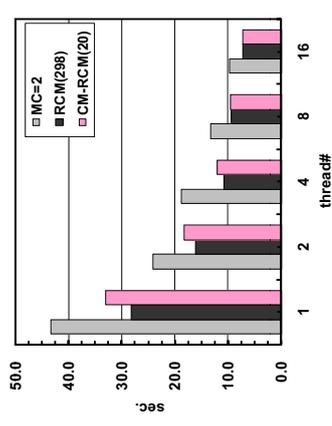
スケラビリティ:numactlの効果

反復回数: MC (2色) : 333回, RCM(298レベル) : 224回
 CM-RCM(Nc=20) : 249回
CASE-1b: 少し良くなった



CASE-0

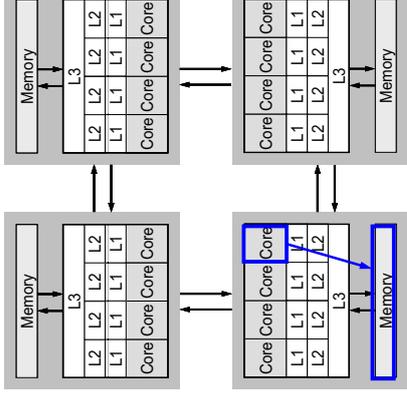
CASE-1b --interleave=all



- L2-solへのOpenMPの実装
- Hitachi SR11000/J2での実行
 - 計算結果
 - CM-RCMオーダリング
- T2K(東大)での実行
- **T2K(東大)での性能向上への道**
 - NUMA Control
 - **First Touch**
 - データ再配置

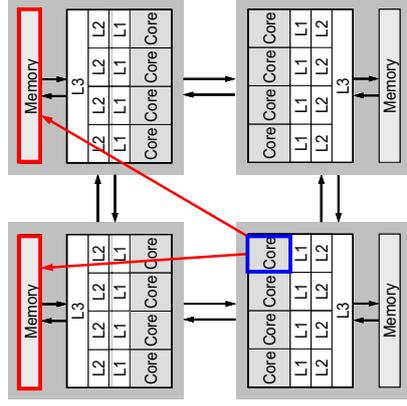
ローカルメモリ, 遠隔メモリ

- コアで扱うデータはなるべくローカルなメモリ(コアの属するソケットのメモリ)上にあると効率が良い。



ローカルメモリ, 遠隔メモリ

- 異なるソケットにある場合はアクセスに時間がかかる。
- First-touchによって、できるだけローカルなメモリ上にデータを持ってくる。



First Touch Data Placement

“Patterns for Parallel Programming” Mattson, T.G. et al.

To reduce memory traffic in the system, it is important to keep the data close to the PEs that will work with the data (e.g. NUMA control).

On NUMA computers, this corresponds to making sure the pages of memory are allocated and “owned” by the PEs that will be working with the data contained in the page.

The most common NUMA page-placement algorithm is the “first touch” algorithm, in which the PE first referencing a region of memory will have the page holding that memory assigned to it.

A very common technique in OpenMP program is to initialize data in parallel using the same loop schedule as will be used later in the computations.

First Touch Data Placement

配列のメモリ・ページ:

最初にtouchしたコアのローカルメモリ上に確保

```
!$omp parallel do private(ip,i,VAL,k)
do ip= 1, PEsmpTOT
  do i = SMPindexG(ip-1)+1, SMPindexG(ip)
    VAL= D(i)*W(i,P)
    do k= indexL(i-1)+1, indexL(i)
      VAL= VAL + AL(k)*W(itemL(k),P)
    enddo
    do k= indexU(i-1)+1, indexU(i)
      VAL= VAL + AU(k)*W(itemU(k),P)
    enddo
    W(i,Q) = VAL
  enddo
enddo
!$omp end parallel do
```

Hokke09

OMP-3

First Touch Data Placement

配列のメモリ・ページ:

最初にtouchしたコアのローカルメモリ上に確保
計算と同じ順番で初期化

```
!$omp parallel do private(ip,i,VAL,k)
do ip= 1, PEsmpTOT
  do i = SMPindexG(ip-1)+1, SMPindexG(ip)
    D(i)=0.d0
    do k= indexL(i-1)+1, indexL(i)
      AL(k) = 0.d0
      itemL(k) = 0
    enddo
    do k= indexU(i-1)+1, indexU(i)
      AU(k) = 0.d0
      itemU(k) = 0
    enddo
  enddo
enddo
!$omp end parallel do
```

Hokke09

98

First Touch

```
!$omp parallel do private(ip,i,VAL,k)
do ip= 1, PEsmpTOT
  do i = SMPindexG(ip-1)+1, SMPindexG(ip)
    VAL= D(i)*W(i,P)
    do k= indexL(i-1)+1, indexL(i)
      VAL= VAL + AL(k)*W(itemL(k),P)
    enddo
    do k= indexU(i-1)+1, indexU(i)
      VAL= VAL + AU(k)*W(itemU(k),P)
    enddo
    W(i,Q) = VAL
  enddo
enddo
!$omp end parallel do
```

OMP-3

99

First Touch

```
!$omp parallel do private(ip,i,VAL,k)
do ip= 1, PEsmpTOT
  do i = SMPindexG(ip-1)+1, SMPindexG(ip)
    VAL= D(i)*W(i,P)
    do k= indexL(i-1)+1, indexL(i)
      VAL= VAL + AL(k)*W(itemL(k),P)
    enddo
    do k= indexU(i-1)+1, indexU(i)
      VAL= VAL + AU(k)*W(itemU(k),P)
    enddo
    W(i,Q) = VAL
  enddo
enddo
!$omp end parallel do
```

• 以下の配列に対する処理が必要

- indexL, indexU, itemL, itemU
- AL, AU, D, BFORCE, PHI (X)
- 以下については既の実施済み
 - W (ICCGの中)

- 青字: ローカルメモリに載ることが保証される変数
- 赤字: ローカルメモリに載ることが保証されない変数
(右辺のp)

First Touch: NFLAG=1 (有り)

```

mn = ICELTOT
allocate (index(0:mn), indexU(0:mn))
indexU(0) = 0
indexU(1) = 1
if (NFLAG.eq.0) then
  do ice=1, ICELTOT
    indexU(ice) = INU(ice)
  enddo
else
  do ice=1, NCOLORTot
    !$omp parallel do private(ip, ice)
    do ip=1, PEsmpTOT
      &
      indexU(ice) = INU(ice)
    enddo
  enddo
endif
do ice=1, ICELTOT
  indexU(ice) = indexU(ice) + indexU(ice-1)
  indexU(ice) = indexU(ice) + indexU(ice-1)
enddo

```

```

NPU= indexU(ICELTOT)
NPU= indexU(ICELTOT)
allocate (item(NPU), AL(NPU))
allocate (PESmp(NPU), D(mn), PHI(mn))
if (NFLAG.eq.0) then
  BFORCE= 0.d0
  D= 0.d0
  PHI= 0.d0
  itemU= 0
  AL= 0.d0
  AU= 0.d0
else
  do ice=1, NCOLORTot
    !$omp parallel do private(ip, ice, k)
    do ip=1, PEsmpTOT
      &
      do ice=1, NCOLORTot
        D(ice)= 0.d0
        PHI(ice)= 0.d0
        BFORCE(ice)= 0.d0
        do k= indexU(ice-1)+1, indexU(ice)
          item(k)= 0
          AL(k)= 0.d0
        enddo
        do ice= indexU(ice-1)+1, indexU(ice)
          item(k)= 0
          AU(k)= 0.d0
        enddo
      enddo
    enddo
  enddo
endif

```

行列ベクトル積の計算法

METHOD=0
solver_ICCG_mc

```

!$omp parallel do private(ip, i, VAL, k)
do ip=1, PEsmpTOT
  do i=1, SMP indexG(ip-1)+1, SMP indexG(ip)
    VAL= D(i)*W(i,P)
    do k= indexL(i-1)+1, indexL(i)
      VAL= VAL + AL(k)*W(itemL(k),P)
    enddo
    do k= indexU(i-1)+1, indexU(i)
      VAL= VAL + AU(k)*W(itemU(k),P)
    enddo
    W(i,0)= VAL
  enddo
enddo
!$omp end parallel do

```

METHOD=1
solver_ICCG_mc_ft

```

do ice=1, NCOLORTot
  !$omp parallel do private(ip, ip1, i, VAL, k)
  do ip=1, PEsmpTOT
    ip1= (ice-1)*PEsmpTOT + ip
    do i= SMP index(ip1-1)+1, SMP index(ip1)
      VAL= D(i)*W(i,P)
      do k= indexL(i-1)+1, indexL(i)
        VAL= VAL + AL(k)*W(itemL(k),P)
      enddo
      do k= indexU(i-1)+1, indexU(i)
        VAL= VAL + AU(k)*W(itemU(k),P)
      enddo
      W(i,0)= VAL
    enddo
  enddo
enddo
!$omp end parallel do
enddo

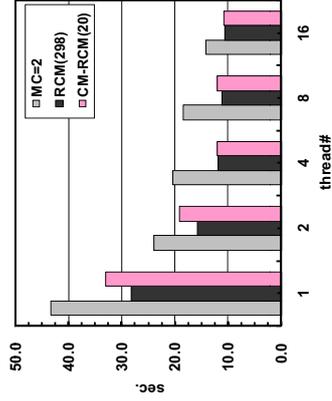
```

スケーラビリティ: First-touchの効果

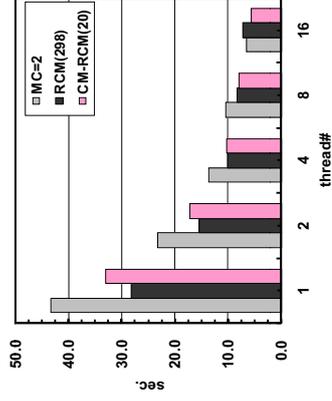
反復回数: MC (2色) : 333回, RCM(298レベル) : 224回

CM-RCM(Nc=20) : 249回

CASE-1a ⇒ 2a: 改善が見られる



CASE-1a
--localalloc



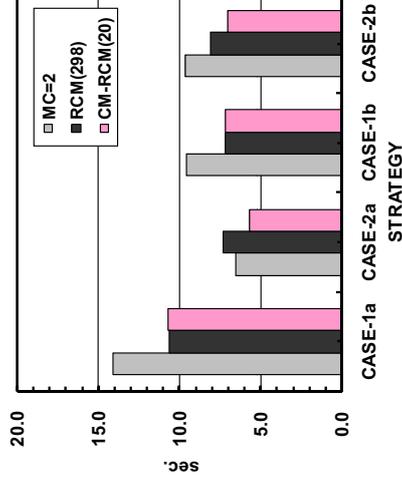
CASE-2a (f5.sh)
METHOD=1の方が良い

スケーラビリティ: First-touchの効果

反復回数: MC (2色) : 333回, RCM(298レベル) : 224回

CM-RCM(Nc=20) : 249回

16コア時における比較, CASE-1b ⇒ 2bは変わらず



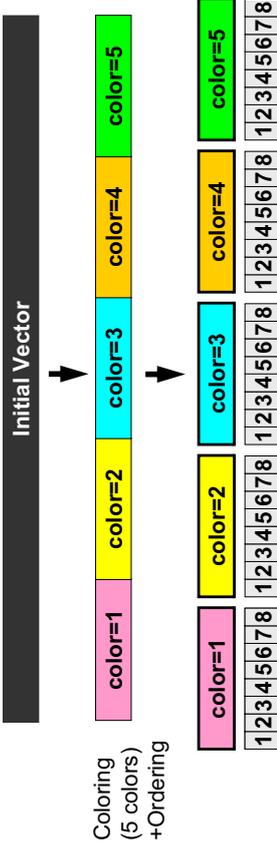
- L2-solへのOpenMPの実装
- Hitachi SR11000/J2での実行
 - 計算結果
 - CM-RCMオーダリング
- T2K(東大)での実行
- **T2K(東大)での性能向上への道**
 - NUMA Control
 - First Touch
 - **データ再配置**

現在のオーダリングの問題

- 色付け
 - MC
 - RCM
 - CM-RCM
- 同じ色に属する要素は独立：並列計算可能
- 「色」の順番に番号付け
- 色内の要素を各スレッドに振り分ける
- 同じスレッド(すなわち同じコア)に属する要素は連続の番号ではない
 - 効率の低下

SMPindex: 前処理向け

```
do ic= 1, NCOLORtot
!$omp parallel do ...
do ip= 1, PEsmpTOT
  ip1= (ic-1)*PEsmpTOT+ip
  do i= SMPindex(ip1-1)+1, SMPindex(ip1)
    (...)
  enddo
enddo
!omp end parallel do
enddo
```



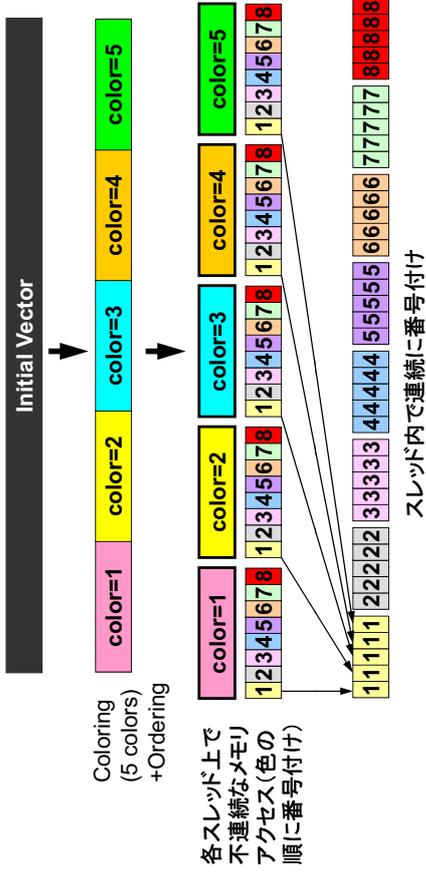
- 5色, 8スレッドの例
- 同じ「色」に属する要素は独立⇒並列計算可能
- 色の順番に並び替え

データ再配置

- 同じスレッドで処理するデータをなるべく連続に配置するように更に並び替え
 - 効率の向上が期待される
- 番号の付け替えによって要素の大小関係は変わるが, 上三角, 下三角の関係は変えない(もとの計算と反復回数は変わらない)
 - 従って自分より要素番号が大きいのにIAL(下三角)に含まれていたりする
- First-touch + データ再配置: CASE3

データ再配置

各スレッド上でメモリアクセスが連続となるよう更に並び替え
5 colors, 8 threads



プログラムのありか

- 所在
 - <\$L3>/reorder
- コンパイル, 実行方法
 - 本体
 - cd <\$L3>/reorder
 - make
 - メッシュ生成
 - <\$L3>/reorder/L3-r-sol (実行形式)
 - コントロールデータ
 - cd <\$L3>/run
 - f90 -O mg.f -o mg
 - 実行用シェル
 - <\$L3>/run/run/INPUT.DAT
 - <\$L3>/run/r5.sh, r6.sh

制御データ(INPUT.DAT)

```

1.00e-00 1.00e-00 1.00e-00 DX/DY/DZ
1.0e-08 EPSICC
16 PEsmpTOT
100 NCOLortot
1 NFLAG
1 METHOD
  
```

変数名	型	内容
DX, DY, DZ	倍精度実数	各要素の3辺の長さ (ΔX, ΔY, ΔZ)
EPSICCG	倍精度実数	収束判定値
PEsmpTOT	整数	データ分割数
NCOLORtot	整数	Ordering手法と色数 ≥ 2 : MC法 (multicolor), 色数 $= 0$: CM法 (Cuthill-Mekee) $= 1$: RCM法 (Reverse Cuthill-Mekee) ≤ -2 : CM-RCM法
NFLAG	整数	0: First-Touch無し, 1: あり
METHOD	整数	行列ベクトル積のループ構造 (0: 従来通り, 1: 前進後退代入と同じ)

再配置

```

allocate (SMP_index(0:PEsmpTOT*NCOLORtot))
do ic=1, NCOLortot
  ni= COLOR_index(ic) - COLOR_index(ic-1)
  num= ni / PEsmpTOT
  do ip=1, PEsmpTOT
    if (.ip. .le. nr) then
      SMP_index((ic-1)*PEsmpTOT+ip)= num + 1
    else
      SMP_index((ic-1)*PEsmpTOT+ip)= num
    endif
  enddo
enddo

allocate (SMP_index_new(0:PEsmpTOT*NCOLORtot))
do ic=1, NCOLortot
  do ip=1, PEsmpTOT
    j0= (ic-1)*PEsmpTOT + ip
    do j=1, ni
      SMP_index_new((ip-1)*NCOLORtot+ip)= SMP_index(ji)
      SMP_index(ji)= SMP_index(j0) + SMP_index(ji)
    enddo
  enddo
  do ip=1, PEsmpTOT
    do ic=1, NCOLortot
      j1= (ip-1)*NCOLORtot + ic
      j0= j1 - 1
      SMP_index_new(j1)= SMP_index_new(j0) + SMP_index(j1)
    enddo
  enddo
  
```

行列ベクトル積の計算法 (METHOD=0)

```

!$omp parallel do private(ip, i, VAL, k)
do ip=1, PEsmpTOT
do i=1, SMPindexG(ip-1)+1, SMPindexG(ip)
VAL= D(i)*W(i, P)
do k= indexL(i-1)+1, indexL(i)
VAL= VAL + AL(k)*W(itemL(k), P)
enddo
do k= indexU(i-1)+1, indexU(i)
VAL= VAL + AU(k)*W(itemU(k), P)
enddo
W(i, 0)= VAL
enddo
!$omp end parallel do

```

Original

```

!$omp parallel do private(ip, i, VAL, k)
do ip=1, PEsmpTOT
do i=1, SMPindex((ip-1)*NCOLORtot)+1, SMPindex(ip*NCOLORtot)
VAL= D(i)*W(i, P)
do k= indexL(i-1)+1, indexL(i)
VAL= VAL + AL(k)*W(itemL(k), P)
enddo
do k= indexU(i-1)+1, indexU(i)
VAL= VAL + AU(k)*W(itemU(k), P)
enddo
W(i, 0)= VAL
enddo
!$omp end parallel do

```

New

前進代入の計算法

```

do ic=1, NCOLORtot
!$omp parallel do private(ip, ip1, i, WVAL, k)
do ip=1, PEsmpTOT
ip1=(ic-1)*PEsmpTOT+ip
do i=1, SMPindex(ip1-1)+1, SMPindex(ip1)
WVAL= W(i, Z)
do k= indexL(i-1)+1, indexL(i)
WVAL= WVAL - AL(k) * W(itemL(k), Z)
enddo
W(i, Z)= WVAL * W(i, DD)
enddo
enddo
!$omp end parallel do
enddo

```

Original

```

do ic=1, NCOLORtot
!$omp parallel do private(ip, ip1, i, WVAL, k)
do ip=1, PEsmpTOT
ip1=(ip-1)*NCOLORtot+ic
do i=1, SMPindex(ip1-1)+1, SMPindex(ip1)
WVAL= W(i, Z)
do k= indexL(i-1)+1, indexL(i)
WVAL= WVAL - AL(k) * W(itemL(k), Z)
enddo
W(i, Z)= WVAL * W(i, DD)
enddo
enddo
!$omp end parallel do
enddo

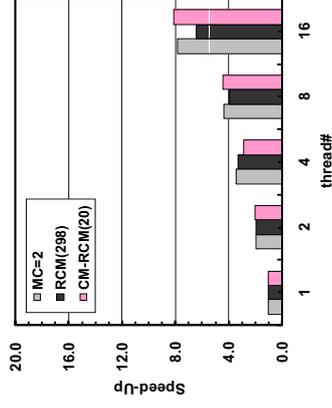
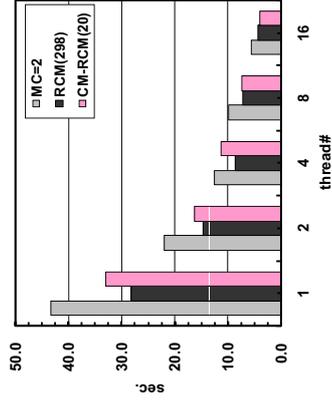
```

New

First-touch+再配置 の効果

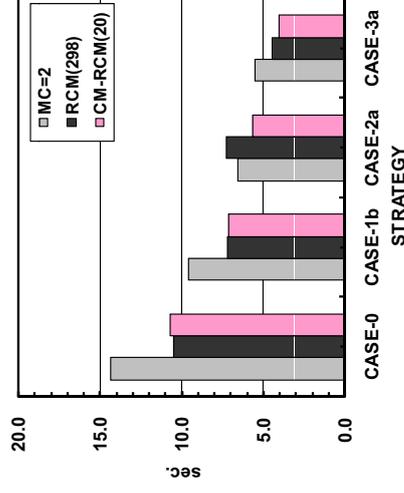
反復回数: MC(2色): 333回, RCM(298レベル): 224回
 CM-RCM(Nc=20): 249回

CASE-3a(--localalloc): METHOD=0の方が良い
 RCMはコア数増えると性能相対的に低下



CASE-0⇒CASE-3:改良の履歴

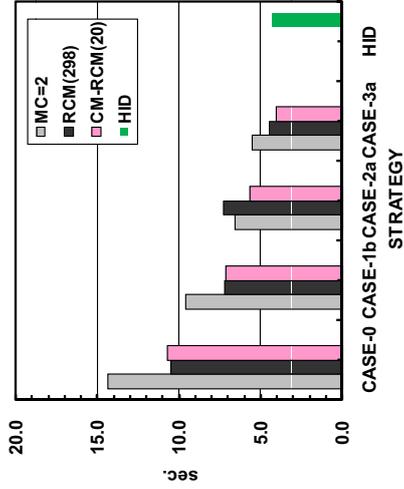
反復回数: MC(2色): 333回, RCM(298レベル): 224回
 CM-RCM(Nc=20): 249回
 16コア時における比較



HID

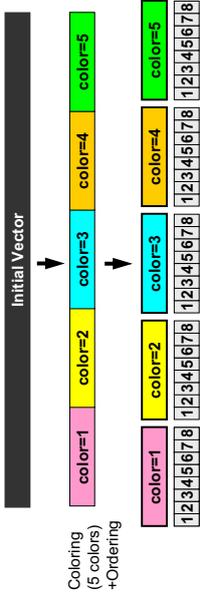
Hierarchical Interface Decomposition

現在取り組んでいる手法: CASE-3aとほぼ同じ
16コア時における比較

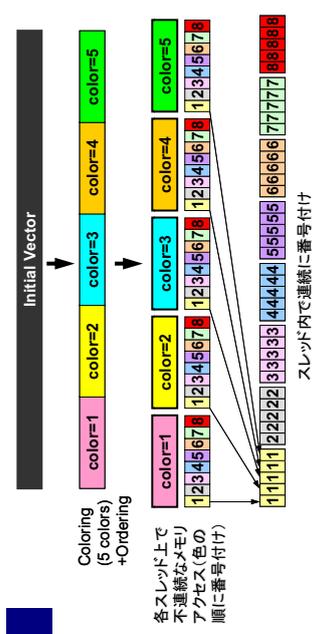


GPUの世界では...

Coalesced
こちらがお勧め

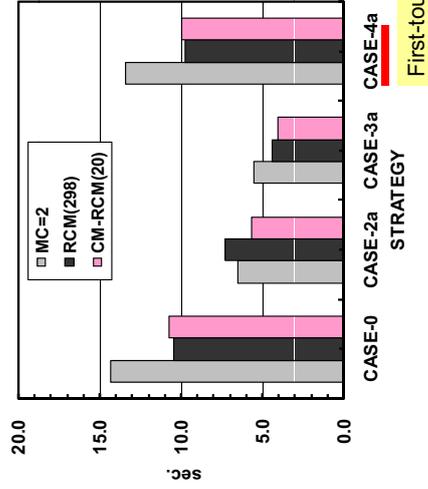


Sequential



First-touch+再配置 の効果

反復回数: MC (2色) : 333回, RCM(298レベル) : 224回
CM-RCM(Nc=20) : 249回
16コア時における比較



- いくつか追加
- STREAM Benchmark
- Hitachi SR1 1000/J2との傾向の比較
- Hopper at Lawrence Berkeley National Laboratoryの結果

STREAM ベンチマーク

<http://www.streambench.org/>

- メモリバンド幅を測定するベンチマーク
 - Copy, Scale, Add, Triad (DAXPYと同じ)

```

Double precision appears to have 16 digits of accuracy
Assuming 8 bytes per DOUBLE PRECISION word
-----
Number of processors = 16
Array size = 2000000
Offset = 0
The total memory requirement is 732.4 MB
( 45.8MB/task)
You are running each test 10 times
-----
The *best* time for each test is used
*EXCLUDING* the first and last iterations
-----
Function      Rate (MB/s)  Avg time  Min time  Max time
Copy:         18334.1898  0.0280    0.0279    0.0280
Scale:        18035.1690  0.0284    0.0284    0.0285
Add:          18649.4455  0.0412    0.0412    0.0413
Triad:        19603.6455  0.0394    0.0392    0.0398

```

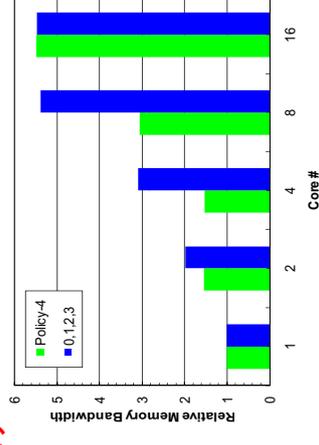
- いくつか追加
- STREAM Benchmark
- Hitachi SR11000/J2との傾向の比較**
- Hopper at Lawrence Berkeley National Laboratoryの結果

Triadの結果 (T2K)

Policy-4@T1の性能(メモリバンド幅)を1.00

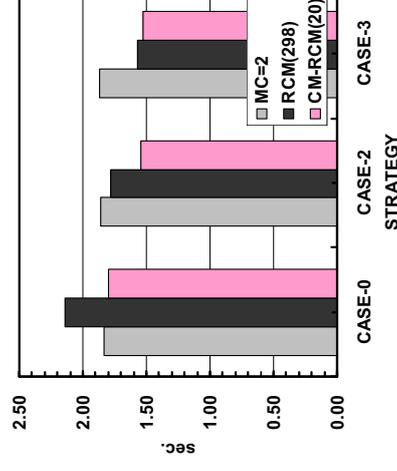
- MPIの場合, NUMA Policyによる比較
- 16コア使った場合, メモリの性能は5倍強程度
- FVMのようなMemory-Boundなアプリケーションは, このメモリ性能に全体の性能が左右される
- CASE-3a: 約8倍@16コア

- 興味のある人は別資料を参照(サンプルプログラムもありますので, 自分でもやってみてください, 後掲ウェブサイトに資料あります)



Hitachi SR11000/J2での結果: 16コア

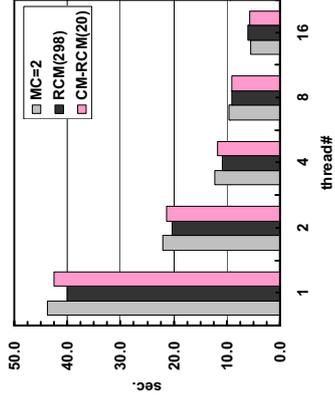
- オーダリング法によって影響が異なる
 - 色数を増やしたときに, First-touch, 再配置の効果が出る



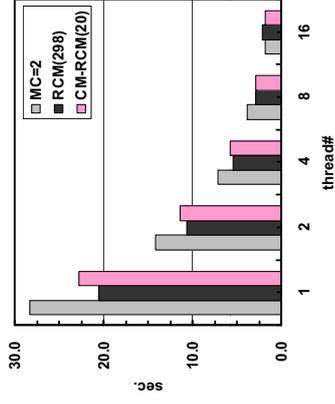
T2KとSR11000 (旧バージョン)

反復回数: MC (2色) : 333回, RCM(298レベル) : 224回
CM-RCM(Nc=20) : 249回

T2K: CASE-3a



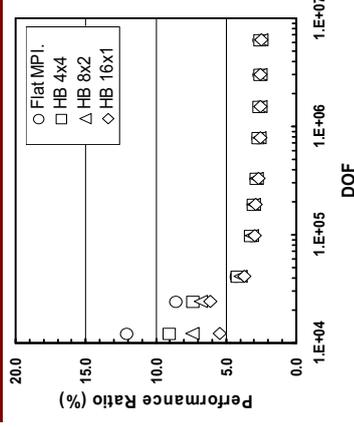
SR11000: CASE-3



GeoFEM Benchmark 1-node with 16-cores: ICCG

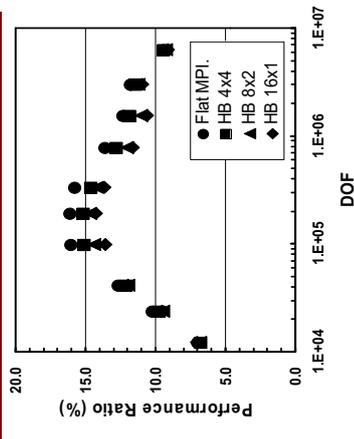
T2K/Tokyo

Opteron 2.3GHz x 16
147.2 GFLOPS/node
20 GB/s for STREAM/Triadd



Hitachi SR11000/J2

Power 5+ 2.3GHz x 16
147.2 GFLOPS/node
100 GB/s for STREAM/Triadd

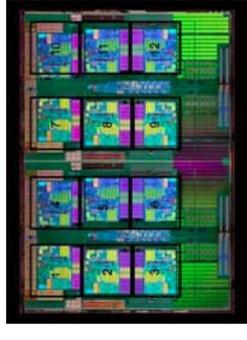


- いくつか追加
- STREAM Benchmark
- Hitachi SR1 1000/J2との傾向の比較
- **Hopper at Lawrence Berkeley National Laboratoryの結果**

Hopper: Cray XE6 at Lawrence Berkeley Natl. Laboratory

<http://www.nersc.gov/systems/hopper-cray-xe6/>

- 6,384 nodes (1.28PFLOPS)
- 2 x 12-core AMD 'MagnyCours' 2.1 GHz processors per node
- 24 cores per node
- 32/64 GB DDR3 1333 MHz

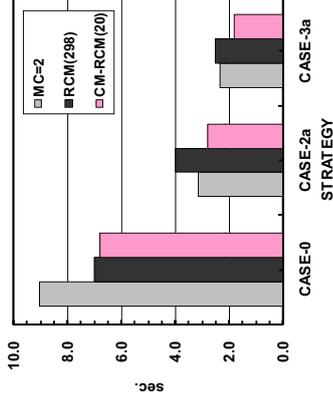


Hopper vs. T2K

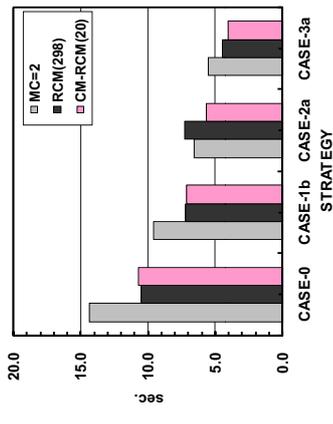
	Hopper	T2K
Peak Performance/core (GFLOPS)	8.4	9.2
Core #/Node	24	16
Peak Performance/node (GFLOPS)	201.6	147.2
Peak Memory		
Bandwidth/node (GB/sec)	85.3	42.7
Triad Performance/node (GB/sec)	52.3	20.0
Triad Performance/core (GB/sec)	2.18	1.25
GeoFEM Bench/ICCG (MFLOPS/core)	469.8	292.8
		1.74
		1.60

Hopper vs. T2K

Hopper with 24-cores

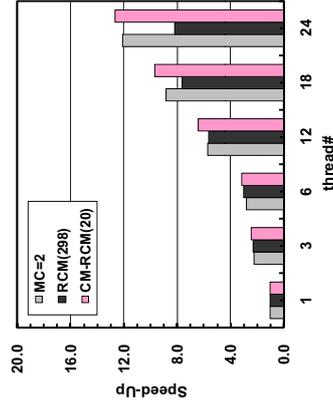


T2K with 16-cores

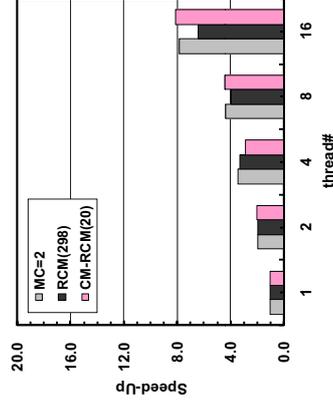


Hopper vs. T2K

Hopper with 24-cores

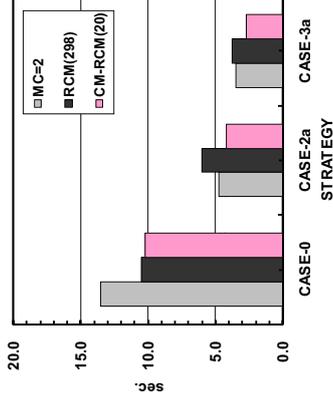


T2K with 16-cores

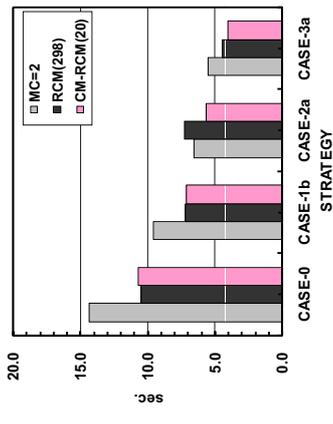


Hopper vs. T2K

Hopper with 16-cores (estimated)



T2K with 16-cores



まとめ

- 「有限体積法から導かれる疎行列を対象としたICCG法」を題材とした、データ配置, reorderingなど, 科学技術計算のためのマルチコアプログラミングにおいて重要なアルゴリズムについての講習
- 更に理解を深めるための, T2Kオーブンスパコン(東大)を利用した実習
- オーダリングの効果
- First-touch
- データ再配置

終わりに

- T2K(東大)
 - 7月8日(金)17:00まで利用可能
- 最終版の資料はWEBにアップ
 - <http://nkl.cc.u-tokyo.ac.jp/seminars/multicore/>
- アンケートにお答えください

今後の動向

- メモリバンド幅と性能のギャップ
 - BYTE/FLOP, 中々縮まらない
- マルチコア化, メニーコア化
- >10⁵コアのシステム
 - Exascale: 10⁸
- **オーダリング**
 - グラフ情報だけでなく, 行列成分の大きさの考慮も必要か?
 - 最適な色数の選択: 研究課題(特に悪条件問題)
- OpenMP+MPIのハイブリッド⇒一つの有力な選択
 - プロセス内(OpenMP)の最適化が最もcritical
- 本講習会の内容が少しでも役に立てば幸いである