

A Heat-Transfer Example with MPI

(short version)

Rolf Rabenseifner
rabenseifner@hls.de

University of Stuttgart
High-Performance Computing-Center Stuttgart (HLRS)
www.hls.de



A Heat-Transfer Example with MPI
Slide 1 Höchstleistungsrechenzentrum Stuttgart



Goals

- first complex MPI example
- most basic MPI features
- base for your own applications
- not a class-room exercise → look on it when traveling home



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 2 Höchstleistungsrechenzentrum Stuttgart



Heat: MPI features

- block data decomposition
- communication: filling the halo with
 - non-blocking point-to-point
 - blocking MPI_SENDRECV
 - MPI_ALLTOALLV
- and for the abort-criterion
 - MPI_ALLREDUCE
 - usage of message tags in point-to-point communication
 - MPI_BCAST
- derived datatypes: MPI_TYPE_VECTOR
- Minor features: MPI_ADDRESS(), MPI_BOTTOM, topologies, MPI_PROC_NULL, MPI_Wtime()
- Performance



Heat Conduction Program - an Example

- solves the partial differential equation for unsteady heat conduction $df/dt = \Delta f$
- uses an explicit scheme: forward-time, centered-space
- solves the equation over a square domain
- initial conditions: $f=0$ everywhere inside the square
- boundary conditions: $f=x$ on all edges
- number of grid points: 80x80
- based on Parallel CFD Test Case
www.hlrs.de/people/resch/PROJECTS/PARACFD.html
- Aims of this session:
 - how to parallelize the example
 - to recapitulate most of the MPI methods



Fortran Introduction for C Programmers

- Array declaration:
 - `type array_name (first_index_dim1 : last_index_dim1, first_index_dim2 : last_index_dim2, ...)`
 - `aaa(17)` means `aaa(1:17)`
- Memory layout: the first index is contiguous
- Using the array:
 - one element: `array_name(i1,i2,...)`
 - sub-array: `array_name (first_sub_index1 : last_sub_index1, first_sub_index2 : last_sub_index2, ...)`
 - referencing:
 - the whole array by the `array_name`
 - a contiguous subarray by the first element `array_name(i1,i2,...)`
- Logical operations: `.gt. / .lt.` means *greater / less than*
`.ge. / .le.` means *greater / less than or equal to*
`.eq. / .ne.` means *equal / not equal to*



Heat: Sequential Program — heat-big.f

```

1      program heat
2      parameter (imax=80,kmax=80)
3      parameter (itmax=20000)
4      double precision eps
5      parameter (eps=1.d-08)
6      double precision
7      phi(0:imax-1,0:kmax-1)
8      dx,dy,dx2,dy2,dx2i,dy2i,dt,
9      dphi,dphimax
10     real tarray(2)
11
12     c
13     dx=1.d0/kmax
14     dy=1.d0/imax
15     dx2=dx*dx
16     dy2=dy*dy
17     dx2i=1.d0/dx2
18     dy2i=1.d0/dy2
19     dt=min(dx2,dy2)/4.d0
20     c start values 0.d0
21     do k=0,kmax-1
22     do i=1,imax-1
23     phi(i,k)=0.d0
24     enddo
25     c start values 1.d0
26     do i=0,imax
27     phi(i,kmax)=1.d0
28     enddo
29     c start values dx
30     phi(0,0)=0.d0
31     do k=1,kmax-1
32     phi(0,k)=phi(0,k-1)+dx
33     enddo
34     c print array
35     write (,'(//,a)')
36     f 'Heat Conduction 2d'
37     write (,'(//,4(a,lpg12.4))')
38     f 'dx =',dx,', dy =',dy,', dt =',dt,',
39     eps =',eps
40     t0=dtime(tarray)
41     c iteration
42     do it=1,itmax
43     dphimax=0.
44     do k=1,kmax-1
45     do i=1,imax-1
46     dphi=(phi(i+1,k)+phi(i-1,k)
47     -2.*phi(i,k))*dy2i
48     +(phi(i,k+1)+phi(i,k-1)
49     -2.*phi(i,k))*dx2i
50     dphi=dphi*dt
51     dphimax=max(dphimax,dphi)
52     phi(i,k)=phi(i,k)+dphi
53     enddo
54     c save values
55     do k=1,kmax-1
56     do i=1,imax-1
57     phi(i,k)=phin(i,k)
58     enddo
59     if(dphimax.lt.eps) goto 10
60     continue
61     t1=dtime(tarray)
62     c print array
63     write (,'(//,a,i6,a)')
64     f 'phi after',it,' iterations'
65     write (,'(//,2(a,lpg12.4,/)')')
66     f 'user time = ',tarray(1),
67     f 'system time = ',tarray(2)
68     stop
69     end

```

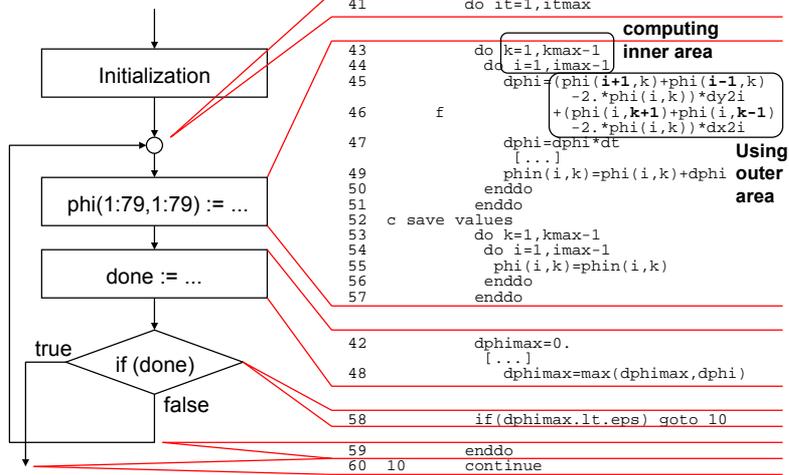
Main loop

Numerical Iteration

Abort criterion



Heat: Sequential Program Scheme

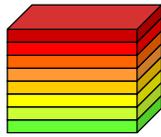


Heat: How to parallelize with MPI

- Block data decomposition of arrays phi, phin
 - array sizes is variable, depends on processor number
- Halo (shadow, ghost cells, ...)
 - to store values of the neighbors
 - width of the halo := needs of the numerical formula (heat: 1)
- Communication:
 - fill the halo with values of the neighbors
 - after each iteration (dt)
- global execution of the abort criterion
 - collective operation

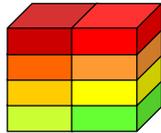


Block data decomposition — in which dimensions?

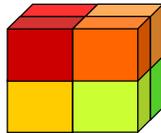


Splitting in

- **one** dimension:
communication
 $= n^2 * 2 * w * 1$



- **two** dimensions:
communication
 $= n^2 * 2 * w * 2 / p^{1/2}$



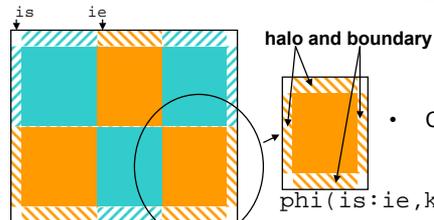
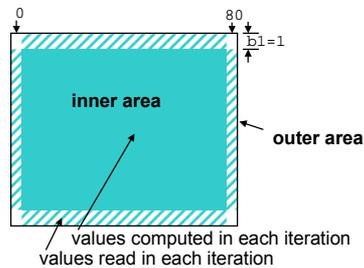
- **three** dimensions:
communication
 $= n^2 * 2 * w * 3 / p^{2/3}$

w=width of halo
n³=size of matrix
p = number of processors
cyclic boundary
→ **two neighbors**
in each direction

optimal for p>11



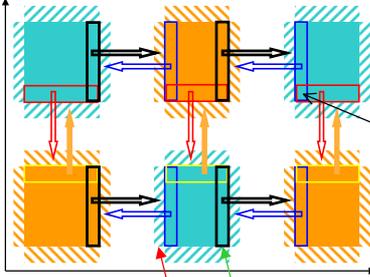
Heat: How to parallelize with MPI - Block Data Decomposit.



- Block data decomposition
- Cartesian process topology
 - two-dimensional
size = idim * kdim
 - process coordinates
icoord = 0 .. idim-1
kcoord = 0 .. kdim-1
 - see heat-mpi-big.f, lines 18-42
- Static array declaration
→ dynamic array allocation
 - phi(0:80,0:80) → phi(is:ie,ks:ke)
 - additional subroutine "algorithm"
 - see heat-mpi-big.f, lines 186+192
- Computation of is, ie, ks, ke
 - see heat-mpi-big.f, lines 44-98



Heat: Communication Method - Non-blocking



→ from left to right
 ← from right to left
 ↓ from upper to lower
 ↑ from lower to upper

- Communicate halo in each iteration
- Ranks via `MPI_CART_SHIFT`
- non-blocking point-to-point
 - Caution:** corners must not be sent simultaneously in two directions
- vertical boundaries:
 - non-contiguous data
 - MPI derived datatypes used

```

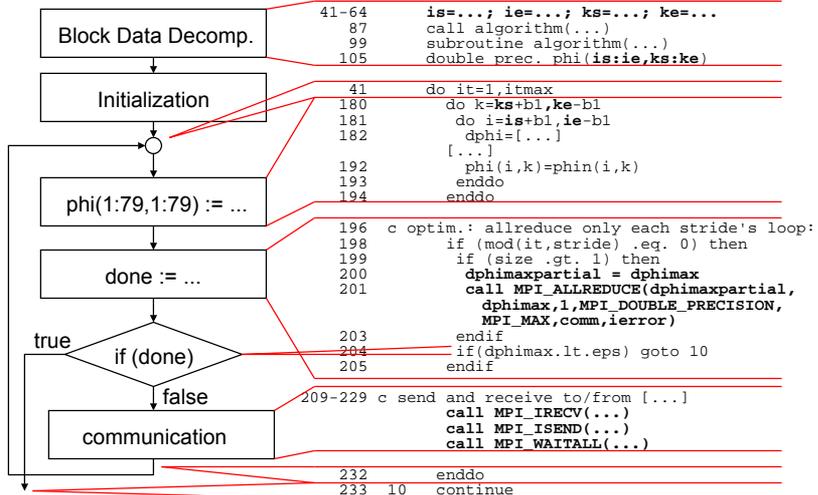
114 call MPI_CART_SHIFT(comm,
    0, 1, left, right, ierror)
115 call MPI_CART_SHIFT(comm,
    1, 1, lower, upper, ierror)
118 call MPI_TYPE_VECTOR(
    kinner, bl, iouter, MPI_DOUBLE_PRECISION,
    vertical_border, ierror)
120 call MPI_TYPE_COMMIT(
    vertical_border, ierror)
221 call MPI_IRECV(phi(is,ks+bl),
    1, vertical_border,
    left, MPI_ANY_TAG, comm, req(1), ierror)
225 call MPI_ISEND(phi(ie-bl,ks+bl),
    1, vertical_border,
    right, 0, comm, req(3), ierror)
  
```



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 11 Höchstleistungsrechenzentrum Stuttgart

HLRS

Heat: Parallel Program Scheme - heat-mpi0-big.f



```

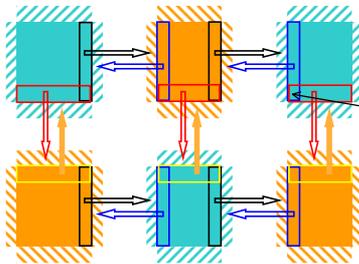
41-64 is=...; ie=...; ks=...; ke=...
87 call algorithm(...)
99 subroutine algorithm(...)
105 double prec. phi(is:ie,ks:ke)
41 do it=1,itmax
180 do k=ks+bl,ke-bl
181 do i=is+bl,ie-bl
182 dphi=[...]
[... ]
192 phi(i,k)=phin(i,k)
193 enddo
194 enddo
196 c optim.: allreduce only each stride's loop:
198 if (mod(it,stride).eq. 0) then
199 if (size.gt. 1) then
200 dphimaxpartial = dphimax
201 call MPI_ALLREDUCE(dphimaxpartial,
    dphimax, 1, MPI_DOUBLE_PRECISION,
    MPI_MAX, comm, ierror)
203 endif
204 if (dphimax.lt.eps) goto 10
205 endif
209-229 c send and receive to/from [...]
call MPI_IRECV(...)
call MPI_ISEND(...)
call MPI_WAITALL(...)
232 enddo
233 10 continue
  
```



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 12 Höchstleistungsrechenzentrum Stuttgart

HLRS

Heat: Different Communication Methods - heat-mpi1-big.f



→ from left to right
 ← from right to left
 ↓ from upper to lower
 ↑ from lower to upper

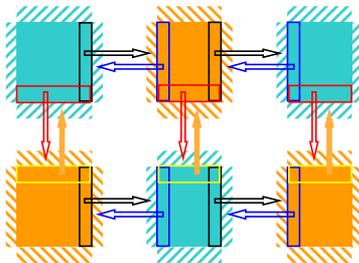
- three different communication methods
 - non-blocking point-to-point
 - examples `communication1-4`
heat-mpi1-big.f, lines 435-442,
lines 523-705
 - **Caution:** corners must not be sent simultaneously in two directions
 - MPI_SENDRECV
 - examples `communication5+6`
heat-mpi1-big.f, lines 435-442,
lines 523-705
 - only for cartesian problems
 - MPI_ALLTOALLV
 - label 207, lines 447-468, 224-265
 - no MPI tag, see next slide
 - Fortran pitfalls with MPI_BOTTOM



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 13 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Heat: sending non-contiguous data - heat-mpi1-big.f



→ from left to right
 ← from right to left
 ↓ from upper to lower
 ↑ from lower to upper

- two methods for vertical boundaries:
- MPI derived datatypes (`comm. 1-5`)
 - some MPI implementations are slow
 - local copying (`comm. 2-7`)
 - non-contiguous array is copied to contiguous scratch array
 - scratch array is sent with MPI
 - may be executed fast
 - optimized by the compiler
 - vectorized
 - automatically multi-threaded parallelized
 - multi-threaded parallelized with OpenMP



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 14 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Heat: Non-blocking point-to-point with derived datatypes (1)

```
214 c create a MPI Vector
215     call MPI_TYPE_VECTOR(kinner,bl,iouter, MPI_DOUBLE_PRECISION,
216                          vertical_border,ierror)
217     call MPI_TYPE_COMMIT(vertical_border, ierror)
218
219     call MPI_TYPE_VECTOR(bl,iinner,iouter, MPI_DOUBLE_PRECISION,
220                          horizontal_border,ierror)
221     call MPI_TYPE_COMMIT(horizontal_border, ierror)

```

```
535     call MPI_IRECV(phi(is+bl,ks),1,horizontal_border,
536                   lower,MPI_ANY_TAG,comm, req(2),ierror)
537     call MPI_IRECV(phi(is+bl,ke),1,horizontal_border,
538                   upper,MPI_ANY_TAG,comm, req(1),ierror)
539     call MPI_ISEND(phi(is+bl,ke-bl),1,horizontal_border,
540                   upper,upper_right_tag,comm, req(3),ierror)
541     call MPI_ISEND(phi(is+bl,ks+bl),1,horizontal_border,
542                   lower,lower_left_tag,comm, req(4),ierror)
543     call MPI_WAITALL(4, req, sts_upper, ierror)

```

```
550     call MPI_IRECV(phi(is,ks+bl),1,vertical_border,
551                   left,MPI_ANY_TAG,comm, req(2),ierror)
552     call MPI_IRECV(phi(ie,ks+bl),1,vertical_border,
553                   right,MPI_ANY_TAG,comm, req(1),ierror)
554     call MPI_ISEND(phi(ie-bl,ks+bl),1,vertical_border,
555                   right,upper_right_tag,comm, req(3),ierror)
556     call MPI_ISEND(phi(is+bl,ks+bl),1,vertical_border,
557                   left,lower_left_tag,comm, req(4),ierror)
558     call MPI_WAITALL(4, req, sts_right, ierror)
```



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 15 Höchstleistungsrechenzentrum Stuttgart



Heat: MPI_SENDRECV with local copying (communica.6)

```
765     double precision phiright(ks+1:ke-1), phircvright(ks+1:ke-1)
766     double precision phileft(ks+1:ke-1), phircvleft(ks+1:ke-1)

```

```
774     call MPI_SENDRECV(phi(is+1,ke-1),iinner,
775                      MPI_DOUBLE_PRECISION,upper,upper_right_tag,
776                      phi(is+1,ks),iinner,
777                      MPI_DOUBLE_PRECISION,lower,MPI_ANY_TAG,
778                      comm,sts_lower,ierror)
779     call MPI_SENDRECV(phi(is+1,ks+1),iinner,
780                      MPI_DOUBLE_PRECISION,lower,lower_left_tag,
781                      phi(is+1,ke),iinner,
782                      MPI_DOUBLE_PRECISION,upper,MPI_ANY_TAG,
783                      comm,sts_upper,ierror)

```

```
790     phileft(ks+1:ke-1)=phi(is+1,ks+1:ke-1)
791     phiright(ks+1:ke-1)=phi(ie-1,ks+1:ke-1)
792     call MPI_SENDRECV(phiright(ks+1),kinner,
793                      MPI_DOUBLE_PRECISION,right,upper_right_tag,
794                      phircvleft(ks+1),kinner,
795                      MPI_DOUBLE_PRECISION,left,MPI_ANY_TAG,
796                      comm,sts_left,ierror)
797     call MPI_SENDRECV(phileft(ks+1),kinner,
798                      MPI_DOUBLE_PRECISION,left,lower_left_tag,
799                      phircvright(ks+1),kinner,
800                      MPI_DOUBLE_PRECISION,right,MPI_ANY_TAG,
801                      comm,sts_right,ierror)
802     if (left.ne.MPI_PROC_NULL)
803         phi(is,ks+1:ke-1)=phircvleft(ks+1:ke-1)
804     if (right.ne.MPI_PROC_NULL)
805         phi(ie,ks+1:ke-1)=phircvright(ks+1:ke-1)
```



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 16 Höchstleistungsrechenzentrum Stuttgart



Heat: MPI_ALLTOALLV with local copying (label 207)

```

448      phileft(ks+1:ke-1)=phi(is+1,ks+1:ke-1)
449      phiright(ks+1:ke-1)=phi(ie-1,ks+1:ke-1)

450 c      ... to prevent compiler optimizations because the compiler
451 c      does not see the input arguments of ALLTOALLV:
452      call MPI_ADDRESS(phi,iaddr_dummy,ierror) <- horizontal inner borders
453      call MPI_ADDRESS(phileft,iaddr_dummy,ierror) <- left vertical inner border
454      call MPI_ADDRESS(phiright,iaddr_dummy,ierror) <- right vertical inner border

455      call MPI_ALLTOALLV(
456      f      MPI_BOTTOM, (sendcounts, sdispls) MPI_DOUBLE_PRECISION,
457      f      MPI_BOTTOM, (recvcounts, rdispls) MPI_DOUBLE_PRECISION,
458      f      comm, ierror)
                                     ↑
                                     arrays: entries for each direction

459 c      ... to prevent compiler optimizations because the compiler
460 c      does not see the output arguments of ALLTOALLV:
461      call MPI_ADDRESS(phi,iaddr_dummy,ierror) <- horizontal halos
462      call MPI_ADDRESS(phirecvleft,iaddr_dummy,ierror) <- left vertical halo
463      call MPI_ADDRESS(phirecvright,iaddr_dummy,ierror) <- right vertical halo

464      if (left.ne.MPI_PROC_NULL)
465      f      phi(is,ks+1:ke-1)=phirecvleft(ks+1:ke-1)
466      if (right.ne.MPI_PROC_NULL)
467      f      phi(ie,ks+1:ke-1)=phirecvright(ks+1:ke-1)

```



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 17 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Heat: Communication Efficiency

Communication-time (number of PEs)	T3E (16)	Hitachi (16)	HP-V (8)
MPI non-blocking	3.4	5.3	2 [sec]
MPI_SENDRECV	1.9	5.8	2 [sec]
MPI_ALLTOALLV	0.8 *)	15.4	2 [sec]
Computation-Time	0.65	1.9	2 [sec]

MPI targets portable and efficient message-passing programming
but

efficiency of MPI application-programming is not portable!

*) up 128 PEs:
ALLTOALLV
is better than
SENDRECV

(measured April 29, 1999, with heat-mpi1-big.f and stride 179,
CRAY T3E: sn6715 hwwt3e.hww.de 2.0.4.48 unicomk CRAY T3E mpt.1.3.0.0.6,
Hitachi: HI-UX/MPP hitachi.rus.uni-stuttgart.de 02-03 0 SR2201,
HP: HP-UX hp-v.hww.de B.11.00 A 9000/800 75859)



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 18 Höchstleistungsrechenzentrum Stuttgart

H L R I S

Heat: Parallel Abort Criterion

- Two methods:
 - MPI_ALLREDUCE
 - exact implementation of the sequential code
 - very expensive
 - less expensive, if not computed in each iteration: using variable `stride`
 - computing efficient `stride` needs knowledge about number of iterations, execution time of MPI_ALLREDUCE, and execution time of one iteration without allreduce:

$$\text{stride}_{\text{optimal}} = \sqrt{\text{iterations} * \text{time_allreduce} / \text{time_exec}}$$
 - Computing **global** abort-criterion by communicating the local abort-criterion via message tags from lowest to highest node and backward



Heat: Parallel Abort Criterion - global & tags

```

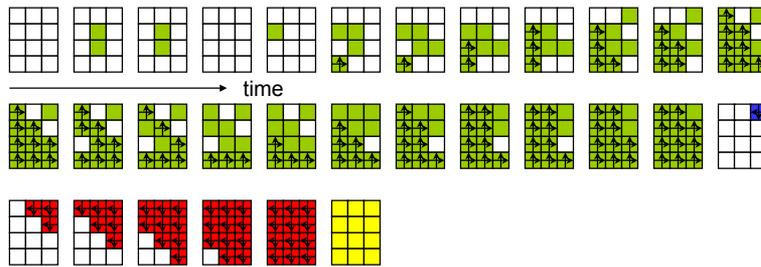
394 103         continue
395 c           ... end_method == 3
396 c           upper_right_tag = 0
397 c           ... outgoing from node (0,0) the local_done criterion
398 c           is propagated with AND operation to the upper-right node
399 c           if ((dphimax.lt.eps) .and.
400 f           ((lower.eq.MPI_PROC_NULL).or.(sts_lower(MPI_TAG).eq.1))
401 f           .and.((left.eq.MPI_PROC_NULL).or.(sts_left(MPI_TAG).eq.1)))
402 f           upper_right_tag = 1
403 c           ... after the local_done criterion has reached at the
404 c           uppermost-rightmost node in the interval of the
405 c           last (idim+kdim) iterations, it is sent back to
406 c           the node (0,0) to inform all nodes about the global stop
407 c           if ((icoord.eq.idim-1) .and. (kcoord.eq.kdim-1)) then
408 f           if ((dphimax.lt.eps) .and.
409 f           ((lower.eq.MPI_PROC_NULL).or.(sts_lower(MPI_TAG).eq.1))
410 f           .and.((left.eq.MPI_PROC_NULL).or.(sts_left(MPI_TAG).eq.1)))
411 f           lower_left_tag = 1
412 c           else
413 f           if (((upper.eq.MPI_PROC_NULL).or.(sts_upper(MPI_TAG).eq.1))
414 f           .and.
415 f           ((right.eq.MPI_PROC_NULL).or.(sts_right(MPI_TAG).eq.1)))
416 f           lower_left_tag = 1
417 c           endif
418 c           it_stopping = it_stopping + lower_left_tag
419 c           ... then the global stop is done synchronously after all
420 c           nodes are informed:
421 c           if (it_stopping .eq. (icoord+kcoord+1)) goto 110
422 c           goto 199

```

1: if ((dphimax.lt.eps) .and. ((lower.eq.MPI_PROC_NULL).or.(sts_lower(MPI_TAG).eq.1)) .and.((left.eq.MPI_PROC_NULL).or.(sts_left(MPI_TAG).eq.1)))
 2: if ((icoord.eq.idim-1) .and. (kcoord.eq.kdim-1)) then if ((dphimax.lt.eps) .and. ((lower.eq.MPI_PROC_NULL).or.(sts_lower(MPI_TAG).eq.1)) .and.((left.eq.MPI_PROC_NULL).or.(sts_left(MPI_TAG).eq.1)))
 3: if (((upper.eq.MPI_PROC_NULL).or.(sts_upper(MPI_TAG).eq.1)) .and. ((right.eq.MPI_PROC_NULL).or.(sts_right(MPI_TAG).eq.1)))
 4: if (it_stopping .eq. (icoord+kcoord+1)) goto 110



Heat: Parallel Abort Criterion - global & tags (cont.)



- Computing,
 - and abort criterion fulfilled locally,
 - and received from left and lower node: Now sending to right and upper node. ①
 - Top-node: Criterion fulfilled locally and received from left and lower node: Now sending stopping-tag to lower and left node. ②
 - Stopping-tag received and now sending to lower and left node. ③
 - Synchronous stopping after all nodes have received the stopping-tag. ④
- At maximum, $2*(dim_1+dim_2-2)$ additional iterations!



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 21 Höchstleistungsrechenzentrum Stuttgart

H L R I S 

Heat: Summary (1)

- block data decomposition
- communication: filling the halo with
 - non-blocking point-to-point
 - blocking MPI_SENDRECV
 - MPI_ALLTOALLV
- and for the abort-criterion and interactive input
 - MPI_ALLREDUCE
 - usage of message tags in point-to-point communication
- derived datatypes
 - MPI_TYPE_VECTOR
 - or copying into contiguous scratch arrays



A Heat-Transfer Example with MPI Rolf Rabenseifner
Slide 22 Höchstleistungsrechenzentrum Stuttgart

H L R I S 

Heat: Summary (2)

- absolute addressing with `MPI_ADDRESS` and `MPI_BOTTOM` and problems with Fortran language binding and MPI (see MPI-2, pages 289f)
- MPI topologies: `MPI_CART_CREATE` / `_COORDS` / `_SHIFT`
- Null processes for not existing neighbors: `MPI_PROC_NULL` (see MPI 1.1, Chapter 3.11, pages 60f)
- timer: `MPI_WTIME`
- performance of different communication methods

