



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

# Praktikum: Paralleles Programmieren für Geowissenschaftler

Prof. Thomas Ludwig, Hermann Lenhart



Dr. Hermann-J. Lenhart

hermann.lenhart@informatik.uni-hamburg.de



## MPI Einführung II:

- Send/Receive Syntax
- Broadcast
- Reduce Operation
- Barrier
- Scatter / Gather -> Matrixbehandlung



## MPI Send/Receive Syntax I

MPI\_SEND(Message, Count, **Datatype**, Dest, Tag, Comm, lerror)

z.B:

Call MPI\_SEND(temp, 1, MPI\_Real, dest, tag, MPI\_COMM\_World, lerror)

temp	Adresse des Sendepuffers; Real :: temp
1	Count – Anzahl der Elemente im Puffer
MPI_Real	Datentyp des gesendeten Elementes
dest	Angabe des Ranges des Zielprozesses; integer :: dest
tag	Nachrichtenkennung; integer :: tag
MPI_COMM_World	Kommunikator (Gruppe, Kontext)
lerror	Fehlerstatus; integer :: lerror



# MPI Send/Receive Syntax II

MPI Datentypen in Anlehnung an Fortran

MPI Datentyp

FORTRAN Datentyp

MPI\_INTEGER

INTEGER

MPI\_REAL

REAL

MPI\_DOUBLE\_PRECISION

DOUBLE PRECISION

MPI\_LOGICAL

LOGICAL

MPI\_CHARACTER

CHARACTER(1)



## MPI Send/Receive Syntax III

MPI\_RECV(Message, Count, Datatype, **Source**, Tag, Comm, **status**, lerror)

Call MPI\_RECV(temp, 1, MPI\_Real, dest, tag, MPI\_COMM\_World, status, lerror)

temp                      Adresse des Sendepuffers;    Real :: Vector

1                              Count – Anzahl der Elemente im Puffer

MPI\_Real                      Datentyp des gesendeten Elementes

**source**                      **Angabe des Ranges des Sendeprozesses; integer :: source**

tag                              Nachrichtenennung (Reihenfolge);    integer :: tag

MPI\_COMM\_World              Kommunikator (Gruppe, Kontext)

**status**                      **Empfangsstatus der Nachricht (angekommen?); integer status(MPI\_STATUS\_SIZE)**

lerror                          Fehlerstatus; integer :: lerror



## MPI Kollektive Operationen

Neben **Point-to-Point Kommunikation** mittels Send & Recv verfügt MPI über umfangreiche Operationen zum **kollektiven Bewegen von Daten**.

**MPI\_BROADCAST** Eine Info an alle Prozesse versenden

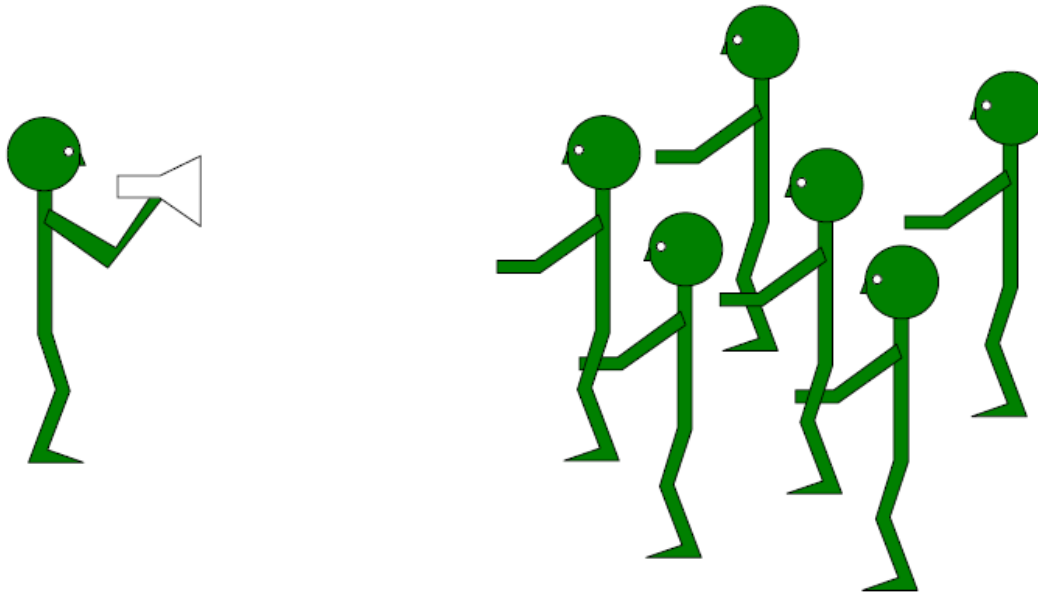
**MPI\_REDUCE** „aggregierende“ Operationen ( Summe; Prod) auf Matrix ausführen

**MPI\_SCATTER** Teilarrays an Prozesse übertragen

**MPI\_GATHER** Teilarrays zusammenführen



# MPI Broadcast I



(Wolfgang Baumann ZIB, 2009;  
Parallel Programming with MPI)



## MPI Broadcast

Neben dem Versenden von Nachrichten zwischen einzelnen Prozesse mittels Call `MPI_SEND(temp, 1, MPI_Real, dest, tag, MPI_COMM_World, lerror)`

gibt es auch die Möglichkeit **eine Nachricht an alle anderen Prozesse**

(z.B. in `MPI_COMM_World`) zu senden:

`MPI_BCAST(Message, Count, Datatype, Root, Comm, lerror)`

Call `MPI_BCAST(temp, 1, MPI_Real, source, MPI_COMM_World, lerror)`

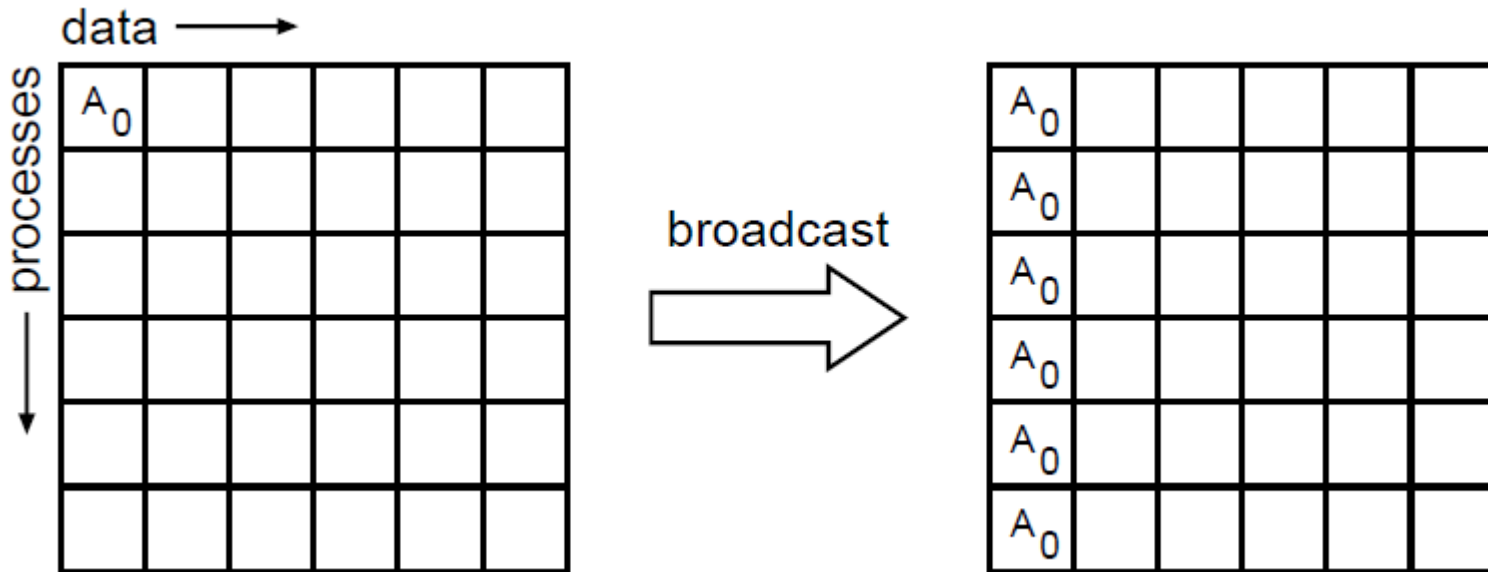
**Für Initialisierung oder zum Programmabbruch genutzt.**





# MPI Broadcast III

Verwendung zur Initialisierung

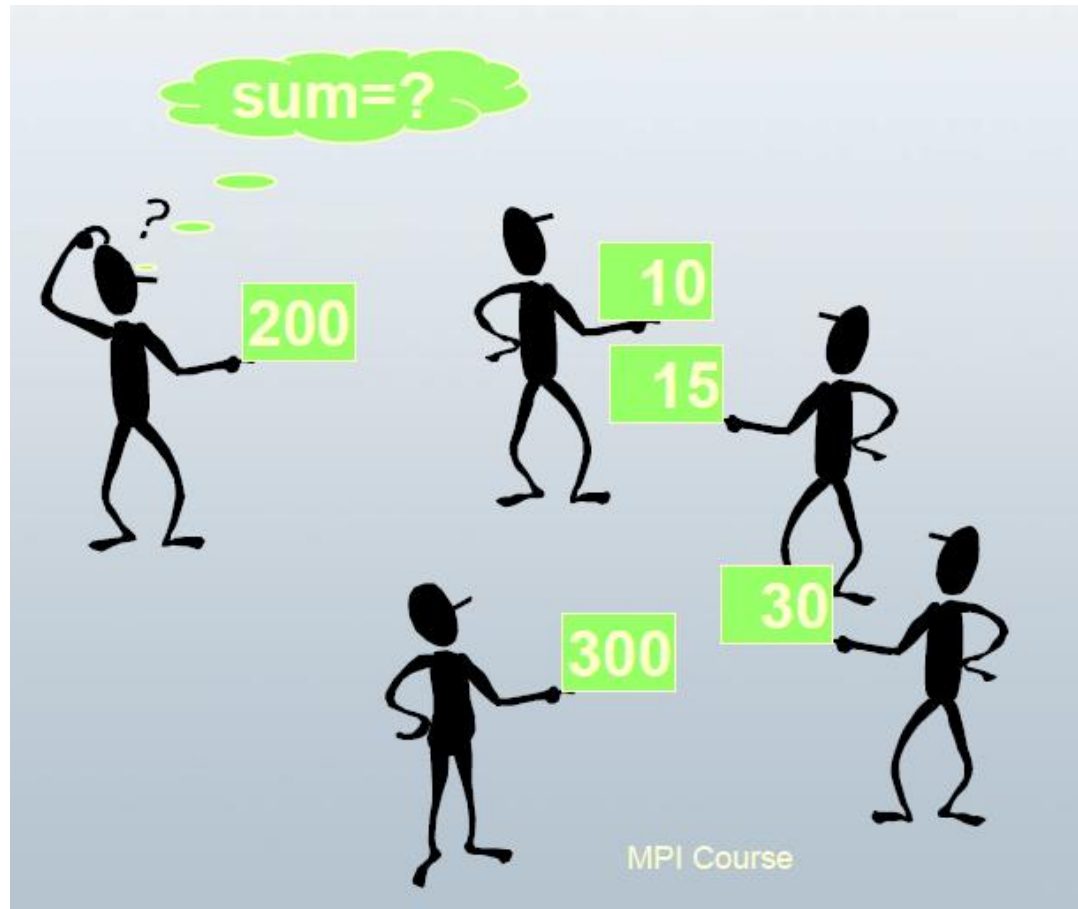


**Es gibt eine interne Hierarchie der Zuteilung!**

(William Gropp ANL, MPI Tutorial)



# MPI Reduce I



DKRZ MPI Einführungs Kurs



## MPI Reduce II

Um die Ergebnisse der einzelnen Prozesse zusammenzuführen gibt es eine Auswahl an „Reduce“ Operationen, z.B:

`MPI_REDUCE(Operand, Result, Count, Datatype, Operation, Root, Comm, Ierror)`

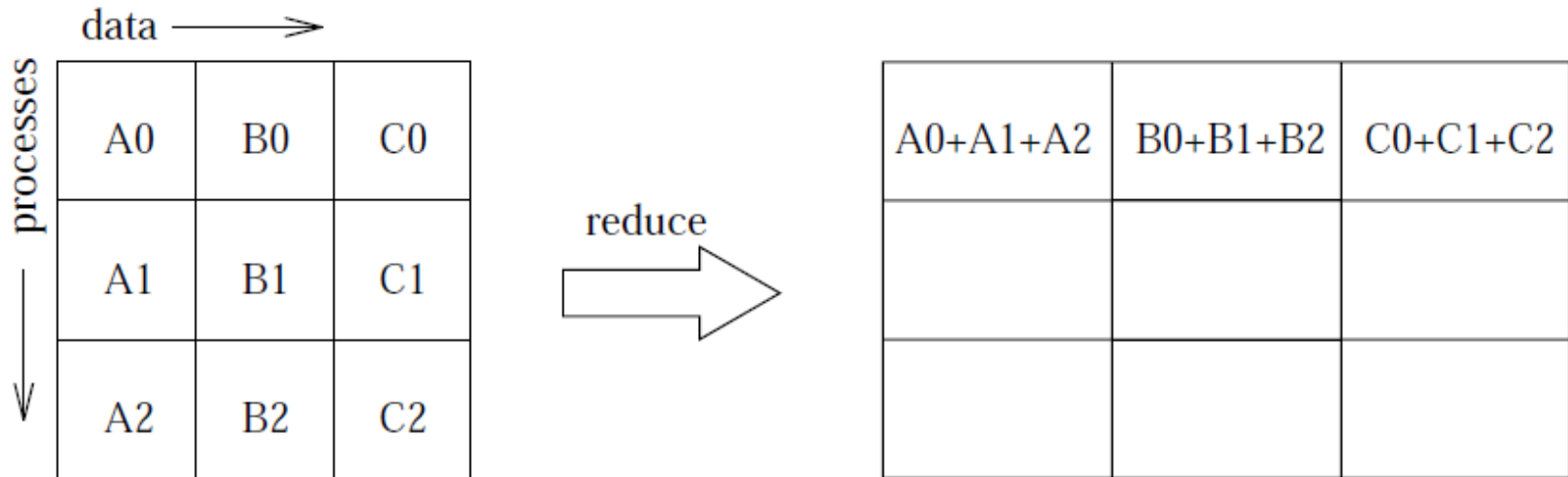
Call `MPI_REDUCE(temp, sum, 1, MPI_Real, MPI_SUM, 0, MPI_COMM_World, Ierror)`

Über die Operation `MPI_SUM` werden alle Resultate der Größe `temp` von allen Prozessen aufaddiert und in der Variable `sum` abgelegt.



# MPI Reduce III

Operator SUM



(William Gropp ANL,  
MPI Tutorial)



## MPI Reduce IV

Übersicht der möglichen MPI\_Reduce Operationen:

MPI_SUM	Summe
MPI_PROD	Produkt
MPI_MAX /MPI_MIN	Maximum/Minimum
MPI_MAXLOC	Maximum und Position des Maximums
MPI_LAND / MPI_LOR	Logical And / Logical Or



## MPI Scatter I

Eine Möglichkeit z.B. eine Anfangsbelegung auf die Teilarrays der Prozesse ,  
zu übertragen bietet MPI\_SCATTER:

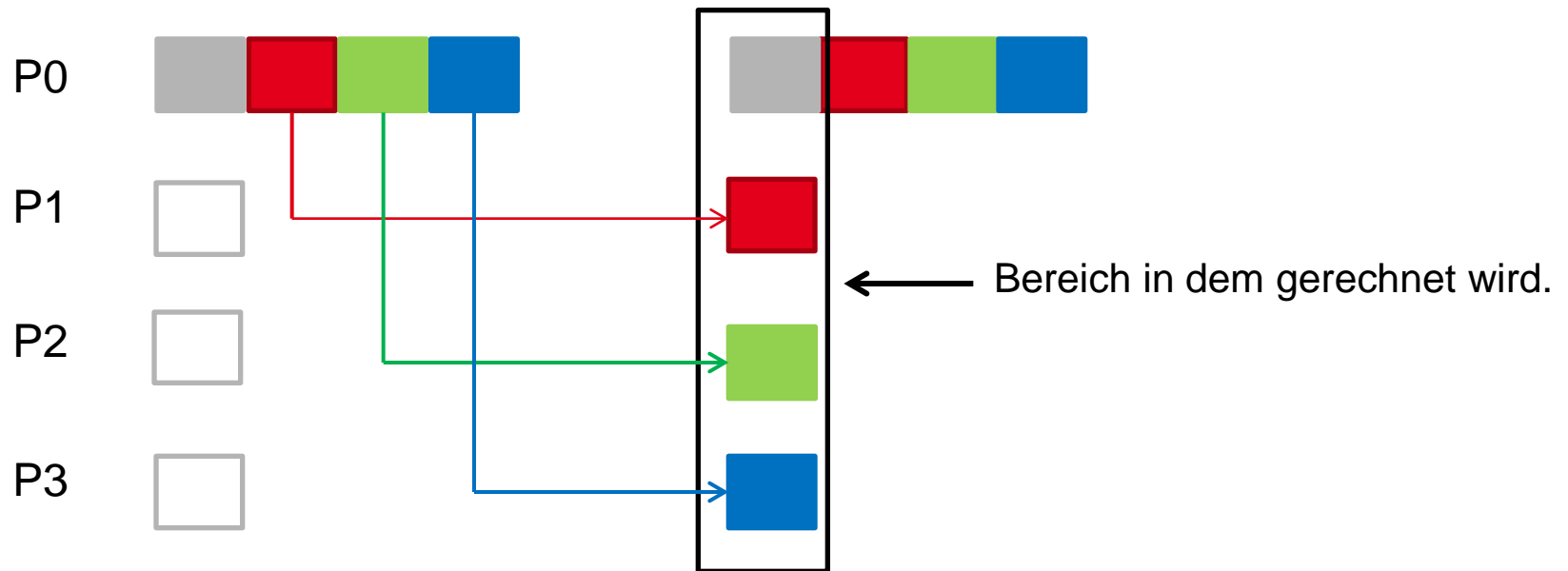
Syntax: MPI\_Scatter(Sendbuffer, Sendcount, Sendtype,  
Recvbuffer, Recvcount, Recvtype,  
Root, Comm, lerror)

Call MPI\_SCATTER(**Send\_Message**, Send\_Count, Send\_Datatype,  
**Recv\_Message**, Recv\_Count, Recv\_Datatype,  
0, MPI\_COMM\_World, lerror)



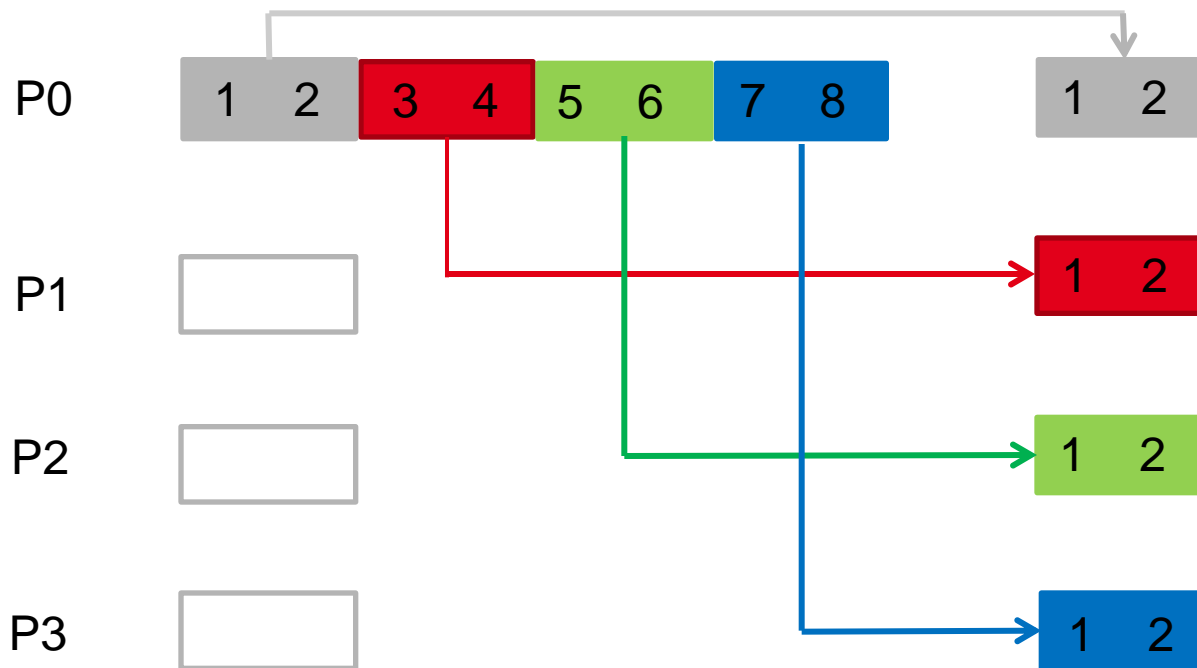
# MPI Scatter II

Die Daten werden von P0 an die anderen Prozesse P1 – P3 gesendet.





# MPI Scatter: array1D(8) => 4 Teile chunk1D(2)





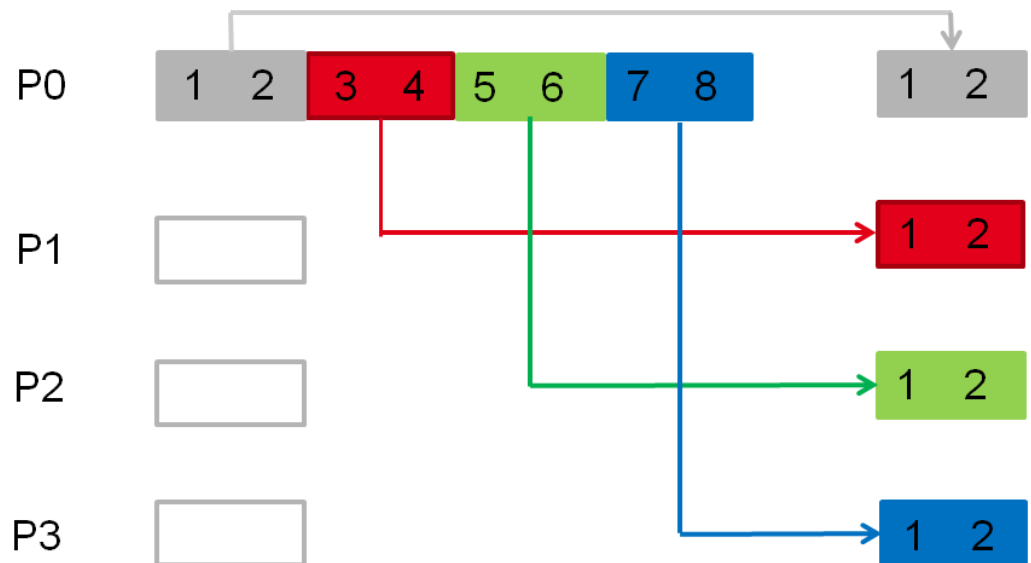
# MPI Scatter:

integer, dimension(8) :: array1D ! all the 1D data  
 integer, dimension(2) :: chunk1D ! piece of 1D data each process works on

! int MPI\_Scatter(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm, ierr)

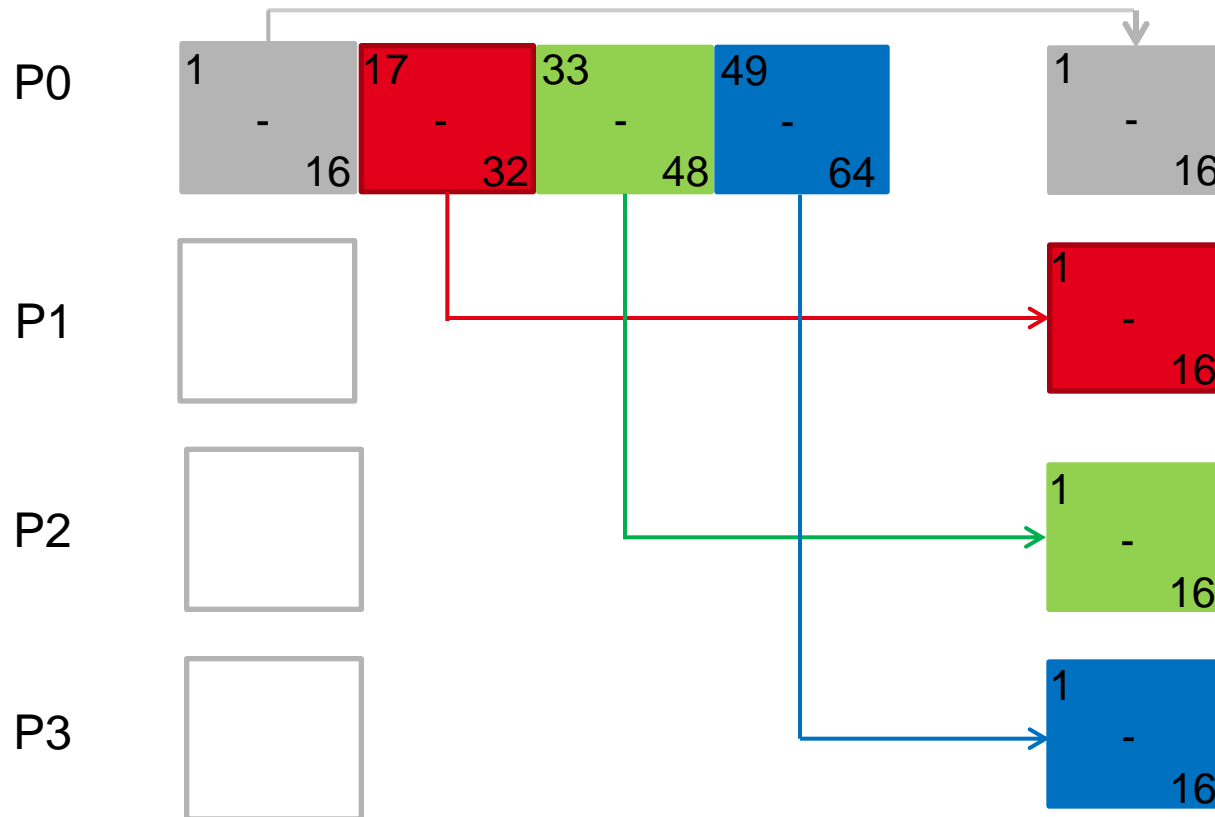
CALL MPI\_SCATTER(array1D, 2, MPI\_INTEGER, chunk1D, 2, MPI\_INTEGER, &0, MPI\_COMM\_WORLD, mpi\_ierr)

Abbildung von Vektor  
array1D auf 4 Teile chunk1D





# MPI Scatter: Beispiel $\text{array2D}(8 \times 8) \Rightarrow 4 \text{ Teile chunk2D}(4 \times 2)$



# MPI Scatter:



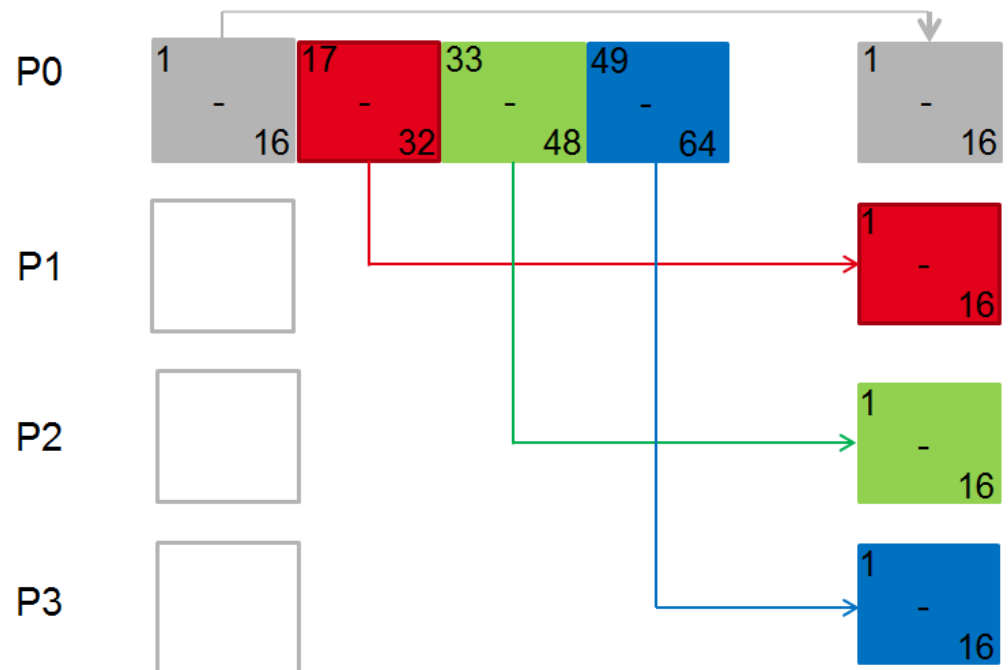
integer, dimension(8, 8) :: array2D ! all the 2D data

integer, dimension(8, 2) :: chunk2D ! piece of 2D data each process works on

! int MPI\_Scatter(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm, ierr)  
 ! 16 = 8\*2 = size(chunk)

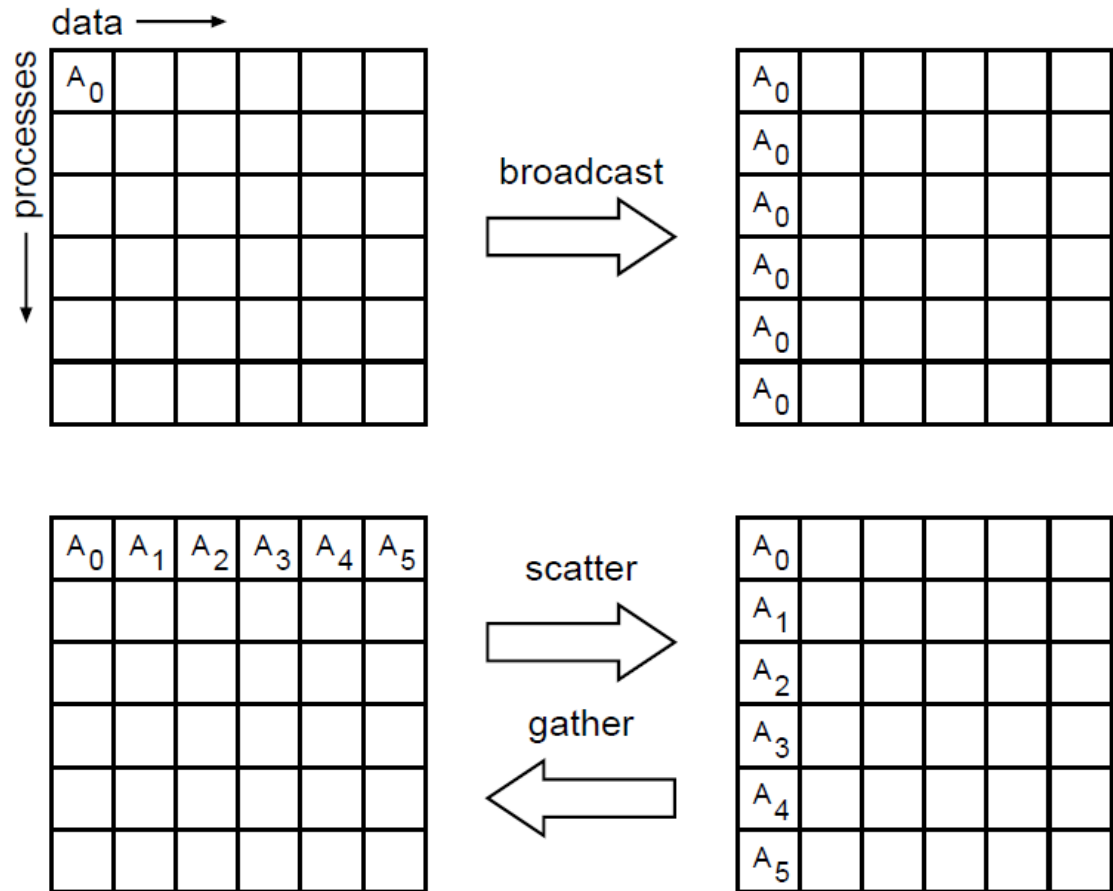
CALL MPI\_SCATTER(array2D, 16, MPI\_INTEGER, chunk2D, 16, MPI\_INTEGER, &  
 0, MPI\_COMM\_WORLD, mpi\_ierr)

Abbildung von Matrix  
 array2D auf 4 Teile chunk2D





# MPI Scatter-Gather



(William Gropp ANL,  
MPI Tutorial)



## MPI Gather I

Eine Möglichkeit Teilarrays der Prozesse (z.B. für I/O Zwecke) wieder zusammenzuführen, bietet MPI\_GATHER:

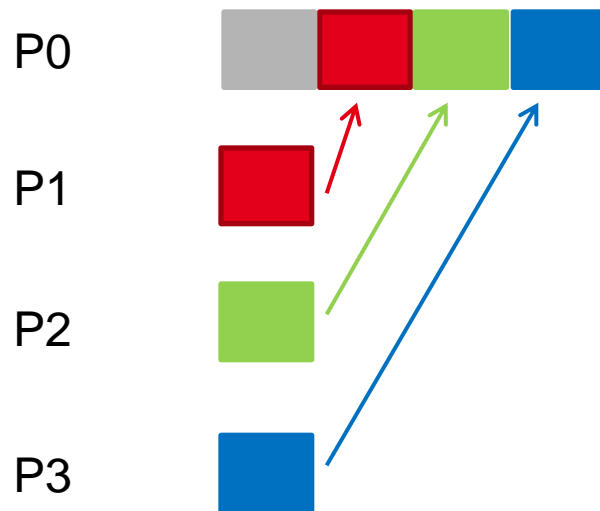
Syntax: MPI\_Gather(**Send\_Message**, Send\_Count, Send\_Datatype,  
**Recv\_Message**, Recv\_Count, Recv\_Datatype,  
 Root, Comm, lerror)

Call MPI\_GATHER (temp, 1, MPI\_Real,  
 temps,1, MPI\_Real,  
 0, MPI\_COMM\_World, lerror)



## MPI Gather II

Call `MPI_GATHER(temp, 1, MPI_Real, tempAll, 1, MPI_Real, 0, MPI_COMM_World, lerror)`





## MPI Barrier I

Der MPI\_BARRIER Befehl wird zur Programmsteuerung eingesetzt.

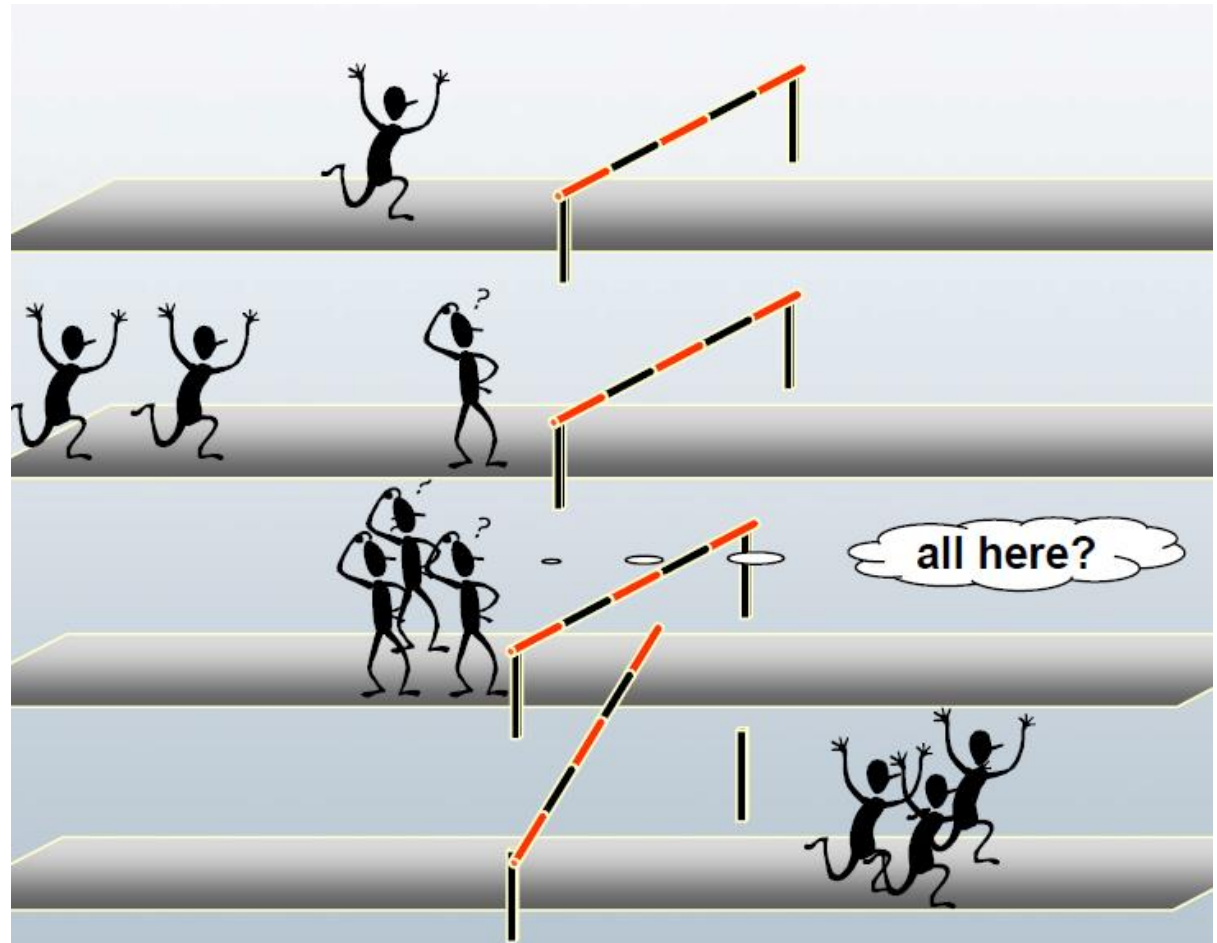
Call MPI\_BARRIER(MPI\_COMM\_World, Ierror)

Der MPI\_Barrier Befehl erzwingt dass alle Prozesse den gleichen Punkt im Code erreicht haben bevor das Programm weiterläuft.



# MPI Barrier II

DKRZ MPI Einführungs Kurs







## MPI Barrier III

Der MPI\_BARRIER Befehl wird vorrangig zur Zeitmessung eingesetzt, z.B.

....

```
Call MPI_BARRIER(MPI_COMM_World, ierror)
```

```
t1 = MPI_WTIME()
```

....

```
Call MPI_BARRIER(MPI_COMM_World, ierror)
```

```
total_time = MPI_WTIME() - t1
```



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG



# Danke das wars!