

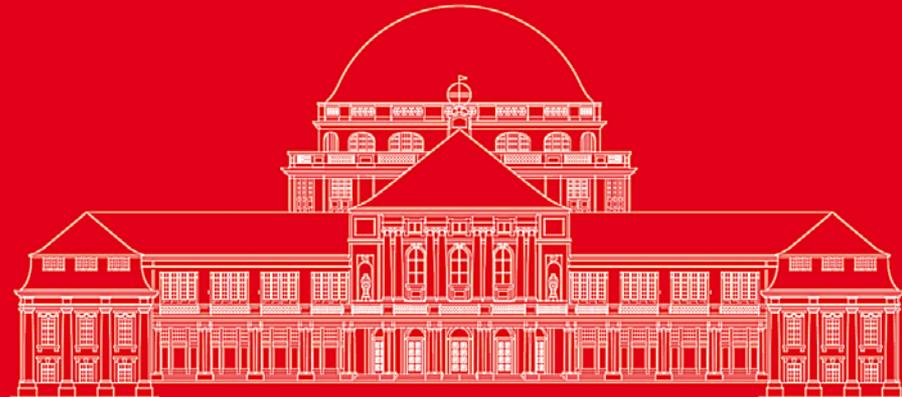


Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

# Praktikum: Paralleles Programmieren für Geowissenschaftler

Prof. Thomas Ludwig, Hermann Lenhart, Ulrich Körner, Nathanael Hübbe



Dr. Hermann-J. Lenhart

[hermann.lenhart@zmaw.de](mailto:hermann.lenhart@zmaw.de)



## MPI Einführung II:

- Kommunikation
- Point to Point Communication
- Nicht blockierende Kommunikation
- MPI Barrier
- Kollektive Operationen



## MPI Kommunikation:

Das wichtigste Kriterium für die Entwicklung paralleler Programme besteht darin die Kommunikation effektiv zu gestalten.

Zur Information:

Ein moderner Parallelrechner kann

bis zu **500 Millionen floating-point Operationen pro Sekunde** berechnen,

aber nur etwa **10 Millionen Wörter pro Sekunde**

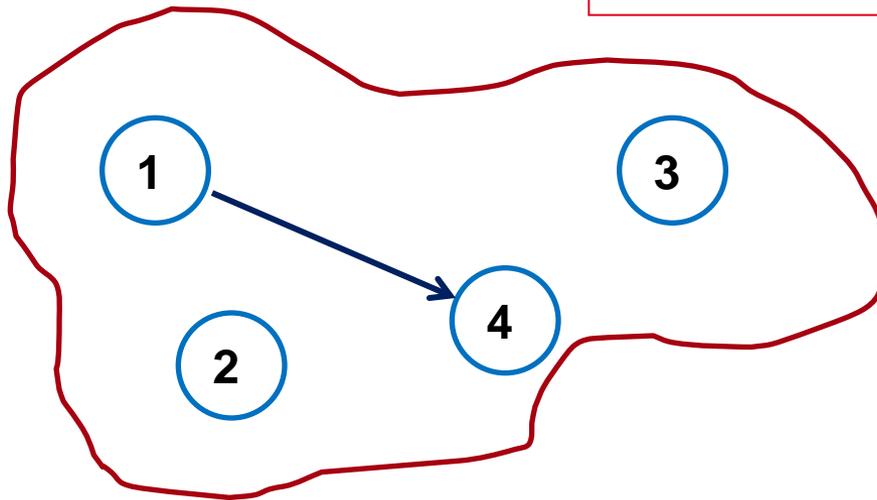
zwischen den Prozessen verschicken!

(Using MPI; Gropp,Lusk,Skjellm, 1999)



# MPI Point to Point Communication:

Kommunikator: MPI\_COMM\_WORLD



Send -> Receive



## MPI Point to Point Communication:

MPI verfügt über **4 verschiedene Send- Funktionen**

Die in unterschiedlichen Send Befehlen dargestellt werden:

Standard

Synchronisiert (synchronous)

Buffered

(Ready)

Alle werden mit dem gleichen Receive Befehl beendet.

„Beendigung“ heißt, der send buffer kann sicher wieder verwendet werden.



## MPI Standard Send 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 Standard Send II

### Beim Standard send ist zu beachten:

Der Standard Send Befehl ist beendet sobald die Message verschickt wurde; unabhängig davon ob die Message schon beim Empfänger angekommen ist, sie kann immer noch für einige Zeit im Netzwerk liegen.



## MPI Synchronous Send I

Wenn sicher gestellt werden soll dass die Message beim empfangenden Prozess angekommen ist, so sollte **synchronous send** verwendet werden.

Wie bei einem „Einschreiben mit Rückantwort“ erhält der Sender eine Bestätigung das die Message angekommen ist.

Erst wenn der Sender diese Bestätigung erhalten hat wird die Sende Operation als abgeschlossen angesehen.

Solange die Bestätigung nicht eingetroffen ist wird der Prozess auf „idle“ gestellt, d.h. das Netzwerk kann niemals mit überfrachtet werden.



## MPI Synchronous Send II

MPI\_**S**SEND(Message, Count, Datatype, Dest, Tag, Comm, lerror)

z.B:

Call MPI\_**S**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 Buffered Send I

Die Option **buffered send** garantiert den Abschluss der Sende Operation sofort, indem die gesendete Message in einen System buffer kopiert wird – für spätere Abwicklung der Zustellung.

Damit ist eine klare Annahme verbunden,  
**Sender und Empfänger können nicht synchron sein.**

Der Programmierer muss dem Programm genug Speicher für den Buffer einräumen:

MPI\_BUFFER\_ATTACH (buffer, *size*, ierror)                      -> Array buffer mit *size* bytes

MPI\_BUFFER\_DETACH (buffer, *size*, ierror)



## MPI Buffered Send II

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

z.B:

Call MPI\_\*\*BSEND(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 Übersicht

Standard Send	MPI_SEND
Synchronous Send	MPI_SSEND
Buffered Send	MPI_BSEND
(Ready Send	MPI_RSEND )
Empfangen	MPI_RECV



## MPI nicht blockierende Kommunikation

Eine weitere Option bietet die nicht-blockierende Kommunikation (non-blocking communication)

Hierbei ist man nicht mehr direkt vom Buffer abhängig, sondern der Programmierer entscheidet selbstständig wann die Nachricht mit dem receive Befehl platziert wird.

Die Syntax dazu bietet: `MPI_ISEND` und `MPI_IRECV`



## MPI non-blocking Send

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

z.B:

Call MPI\_ISEND(temp, 1, MPI\_Real, dest, tag, MPI\_COMM\_World, req, 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)
<b>request</b>	<b>Handle; integer :: req</b> , lerror



## MPI non-blocking Receive

MPI\_IRECV(Message, Count, Datatype, Source, Tag, Comm, **request**, Ierror)

z.B:

Call MPI\_IRECV(temp, 1, MPI\_Real, source, tag, MPI\_COMM\_World, **req**, Ierror)

temp	Adresse des Sendepuffers; Real :: temp
1	Count – Anzahl der Elemente im Puffer
MPI_Real	Datentyp des gesendeten Elementes
source	Angabe des Ranges des Sourceprozesses; integer :: source
tag	Nachrichtenkennung; integer :: tag
MPI_COMM_World	Kommunikator (Gruppe, Kontext)
<b>request</b>	<b>Handle; integer :: req !!! Kein Status</b>
	, Ierror



## MPI non-blocking Kommunikation Statusabfrage

Zur Abfrage des **Status der Isend Nachricht** , gibt es folgende Option:

MPI\_WAIT(request,status, lerror)

Integer :: request, status(MPI\_STATUS\_SIZE), ierror

CALL MPI\_ISEND(buffer, count, datatype, dest, tag, comm, request, ierr)

..... **Der Prozess rechnet weiter**

Call WAIT ( request, status, ierr)



Geht weiter wenn die ISEND-Nachricht angekommen ist  
und der Buffer wieder frei ist und neu belegt werden kann.



## MPI non-blocking Kommunikation Statusabfrage

Zur Abfrage des Status von mehreren Nachricht , gibt es **MPI\_WAITALL**:

Integer :: nx, s, e, comm1d, nbrtop, nbrbot, status\_array(MPI\_STATUS\_SIZE,4), req(4)

Real: a(0:nx+1,s-1:e+1)

Call MPI\_Irecv (a(1,s-1), nx, MPI\_REAL,nbrbot,0,comm1d,req(1),ierr)

Call MPI\_Irecv (a(1,e+1), nx, MPI\_REAL,nbrtop,1,comm1d,req(2),ierr)

Call MPI\_Send(a(1,e), nx, MPI\_REAL, nbrtop, 0, comm1d, req(3), ierr)

Call MPI\_Send(a(1,s), nx, MPI\_REAL, nbrbot, 1, comm1d, req(4), ierr)

Call MPI\_WAITALL(4,req,status\_array, ierr)



## 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

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
```



## MPI Kollektive Operationen

MPI verfügt über umfangreiche Operationen zum kollektiven Bewegen von Daten.

Dazu gehören:

MPI\_BROADCAST

MPI\_REDUCE

MPI\_GATHER

MPI\_SCATTER



## MPI Gather I

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

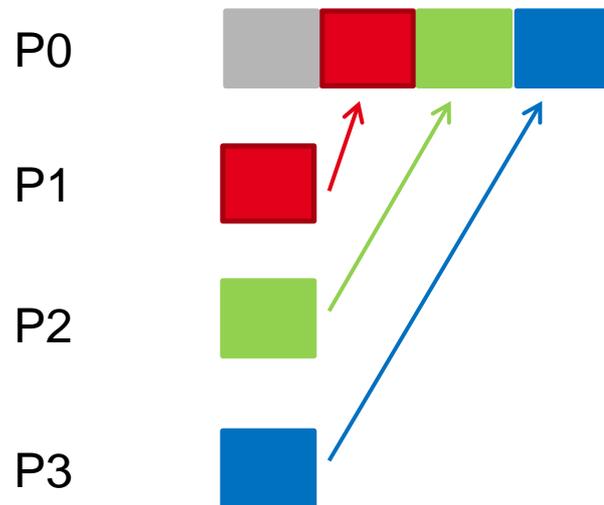
Syntax: MPI\_Gather(Sendmessage, Sendcount, Sendtype,  
Recvmessage, Recvcount, Recvtype,  
Root, Comm, lerror)

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 Scatter I

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

Syntax: MPI\_Scatter(Sendmessage, Sendcount, Sendtype,  
Recvmessage, Recvcount, Recvtype,  
Root, Comm, Ierror)



## MPI Scatter II

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

