

# **4.4. MPI**

## Message Passing Interface

Ferienakademie 2009

Franz Diebold

---

# Agenda

1. Einführung, Motivation
2. Kommunikationsmodell
3. *Punkt-Zu-Punkt*-Kommunikation
4. Globale Kommunikation
5. Vergleich MPI und OpenMP

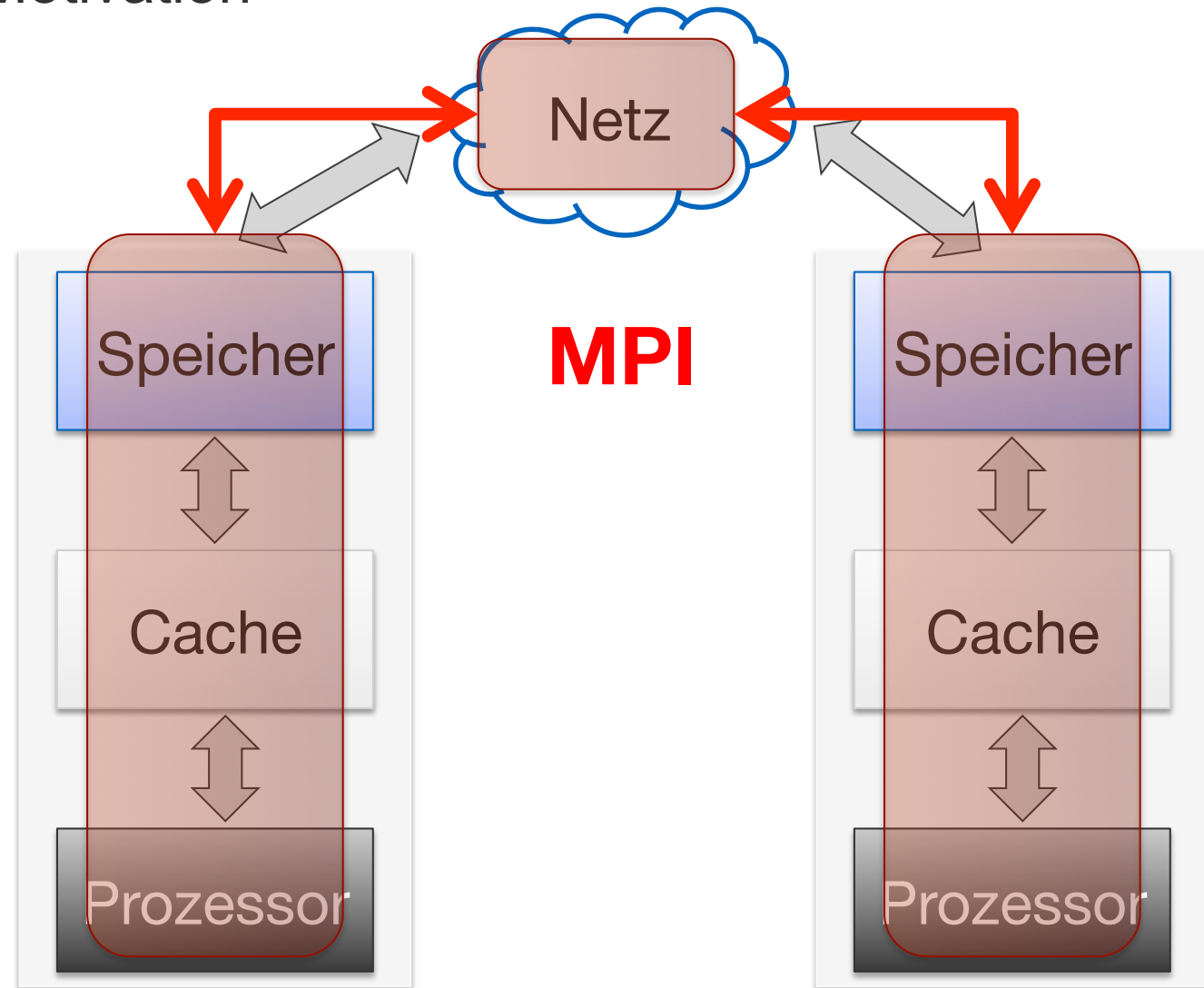
# 1. Einführung, Motivation

## MPI = Message Passing Interface:

- *De-facto-Standard* (Spezifikation) für Nachrichtenaustausch auf verteilten Computersystemen
- Sowohl auf *verteiltern*, als auch auf *gemeinsamem Speicher* nutzbar
- Ziele:
  - Einfache Anwendbarkeit
  - Portabilität
  - Flexibilität
  - Effizienz
- MPI 1 (1994), MPI 2 (1997)
- Implementierungen: C/C++, Fortran, Ocaml, Java, Python

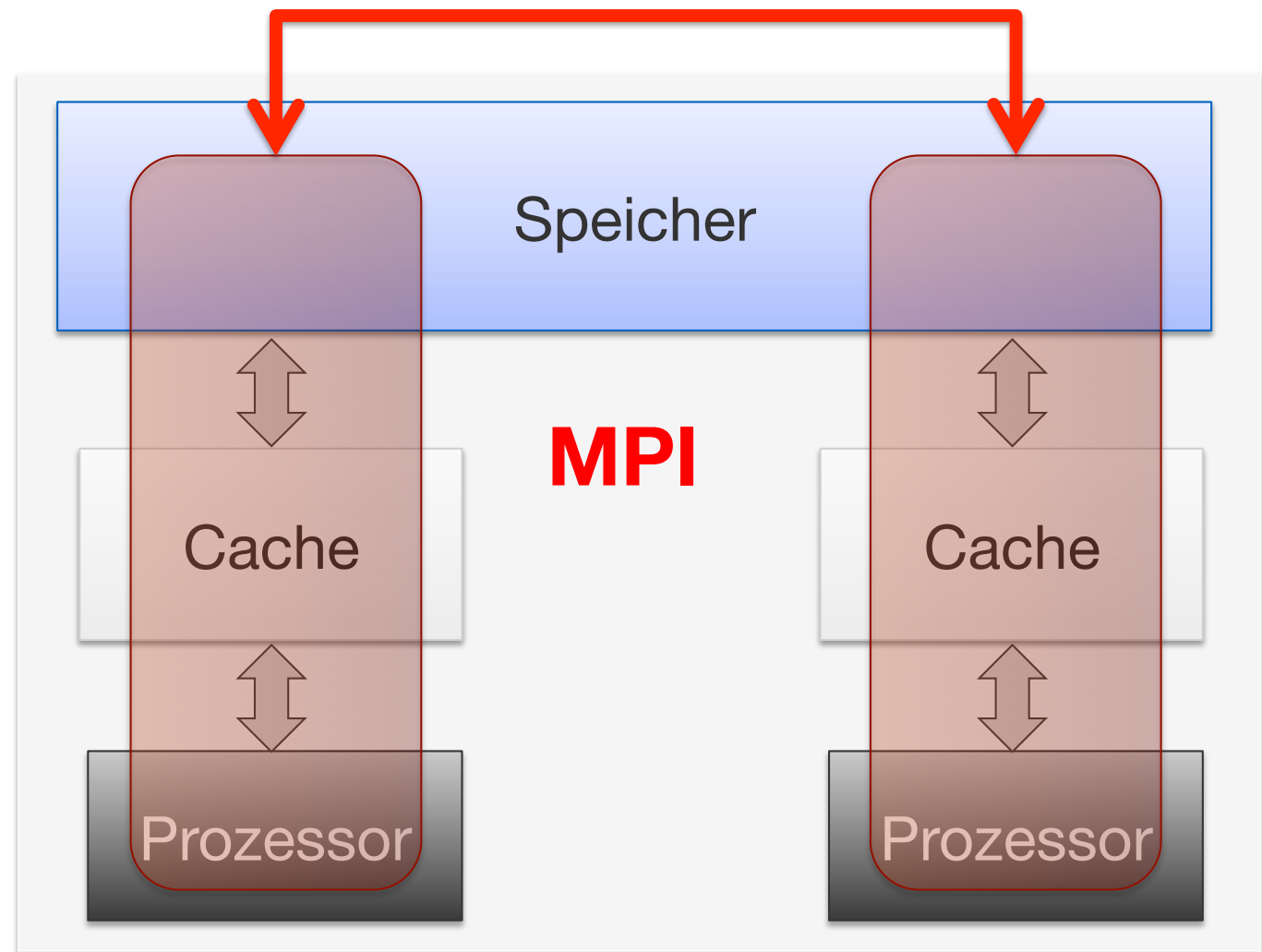
# 1. Einführung, Motivation

Verteilter Speicher,  
über ein Netz  
verbunden



# 1. Einführung, Motivation

Gemeinsamer  
Speicher:  
z.B. SMP  
(Symmetrisches  
Multiprozessor-  
system)



# 1. Einführung, Motivation

## Was MPI bietet:

- *Schnittstelle*
- Regelung des Nachrichtenverkehrs
- Verschiedene *Kommunikationsverfahren*

## Was MPI **nicht** bietet:

- *Automatische Parallelisierung*
- Sicherstellung der Korrektheit der Parallelität
- Automatische *Deadlock-Erkennung & -Vermeidung*
- Automatische *Datenverteilung* über die beteiligten Prozesse

# 1. Einführung, Motivation

## Paradigma von MPI-Programmen: **SPMD**

- Single Programm/Process Multiple Data
- Verzweigung innerhalb des Programmes, um verschiedene Rechner/Knoten unterscheiden zu können

323 Funktionen in Version 2:

- MPI-Funktionen in C:

**MPI\_Name** (Parameter [ , ... ])

# Agenda

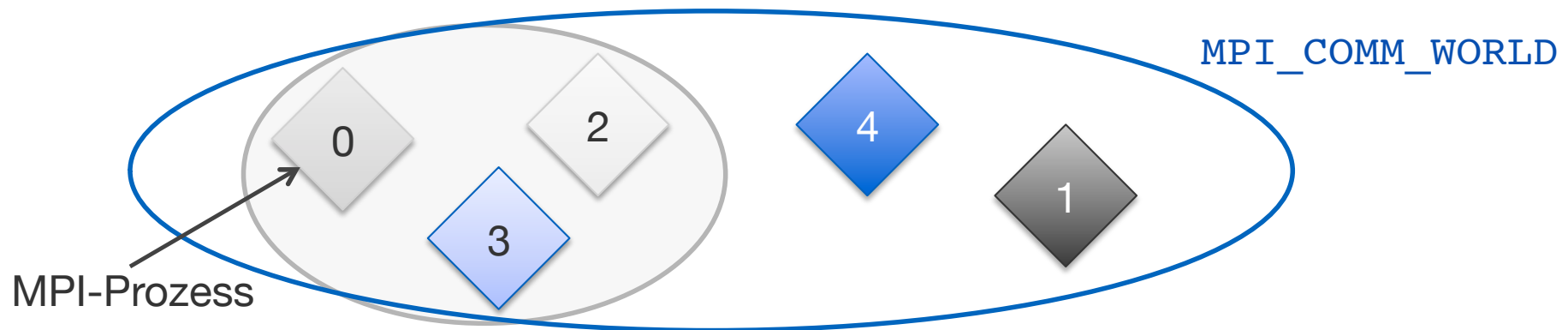
1. Einführung, Motivation
2. Kommunikationsmodell
3. *Punkt-Zu-Punkt*-Kommunikation
4. Globale Kommunikation
5. Ausblick



## 2. Kommunikationsmodell: Begriffe

### Kommunikator:

- Menge von MPI-Prozessen
- Definiert *Kontext* für Kommunikation
- Vordefinierte Kommunikatoren:
  - `MPI_COMM_WORLD` → alle gestartete Prozesse
- Ermittlung der Größe eines Kommunikators im Programm:  
`MPI_Comm_size` (comm, &size)

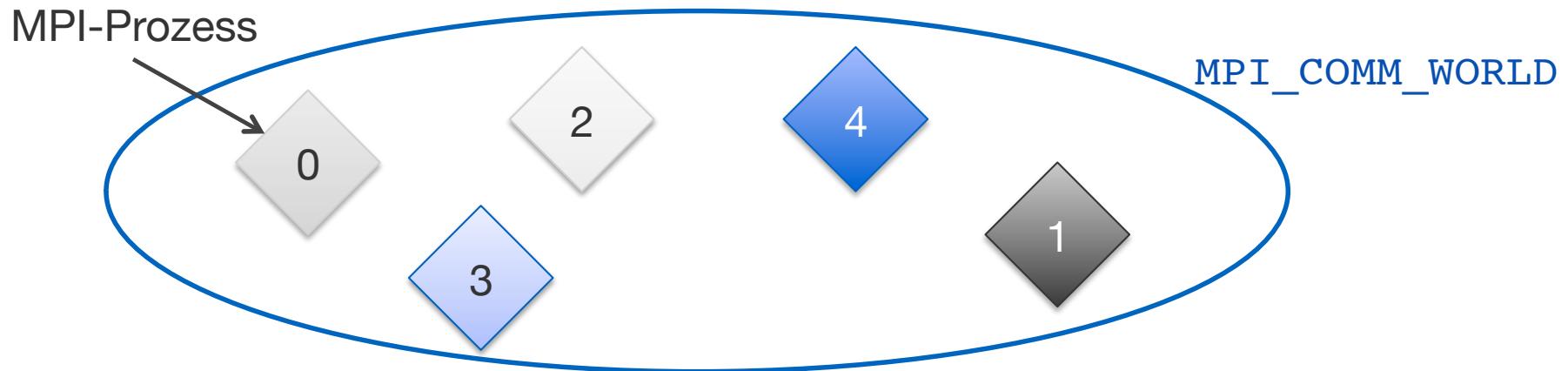


## 2. Kommunikationsmodell: Begriffe

### Rang:

- *Eindeutige Kennung* eines Prozesses
- Bei Initialisierung zugewiesen
- Wird als Quelle und Ziel einer Nachricht angegeben
- Ermittlung im Programm:

**`MPI_Comm_rank`** (`comm`, `&rank`)



## 2. Kommunikationsmodell

### Grundlegende MPI-Funktionen (in C):

- Initialisierung der MPI-Umgebung: (z.B. Verbindungsaufbau)

```
MPI_Init (&argc, &argv);
```

- Muss erste MPI-Funktion im Programm sein
- Aufrufparameter werden an MPI-Laufzeitumgebung übergeben

- Beenden der MPI-Umgebung: (z.B. Verbindungsabbau)

```
MPI_Finalize ();
```

- Muss letzte MPI-Funktion im Programm sein

## 2. Kommunikationsmodell

### Grundlegender Aufbau eines MPI-Programms in C:

```
#include <mpi.h>           // Header-Datei einbinden

// ...

MPI_Init (&argc, &argv);  // MPI initialisieren

// ...
// „MPI-Programm“ hier ...
// ...

MPI_Finalize ();           // MPI beenden
```

## 2. Kommunikationsmodell

### Einfaches MPI-Programm in C:

```
#include <mpi.h>
#include <stdio.h>

int main(int argc, char *argv[]) {
    int numtasks, rank;
    MPI_Init (&argc, &argv);
    MPI_Comm_size (MPI_COMM_WORLD, &numtasks);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    printf ("Hallo Welt! Ich bin Prozess %d von %d.\n", rank,
numtasks);
    MPI_Finalize ();
}
```

## 2. Kommunikationsmodell

### Einfaches MPI-Programm in C: Beispielausgabe

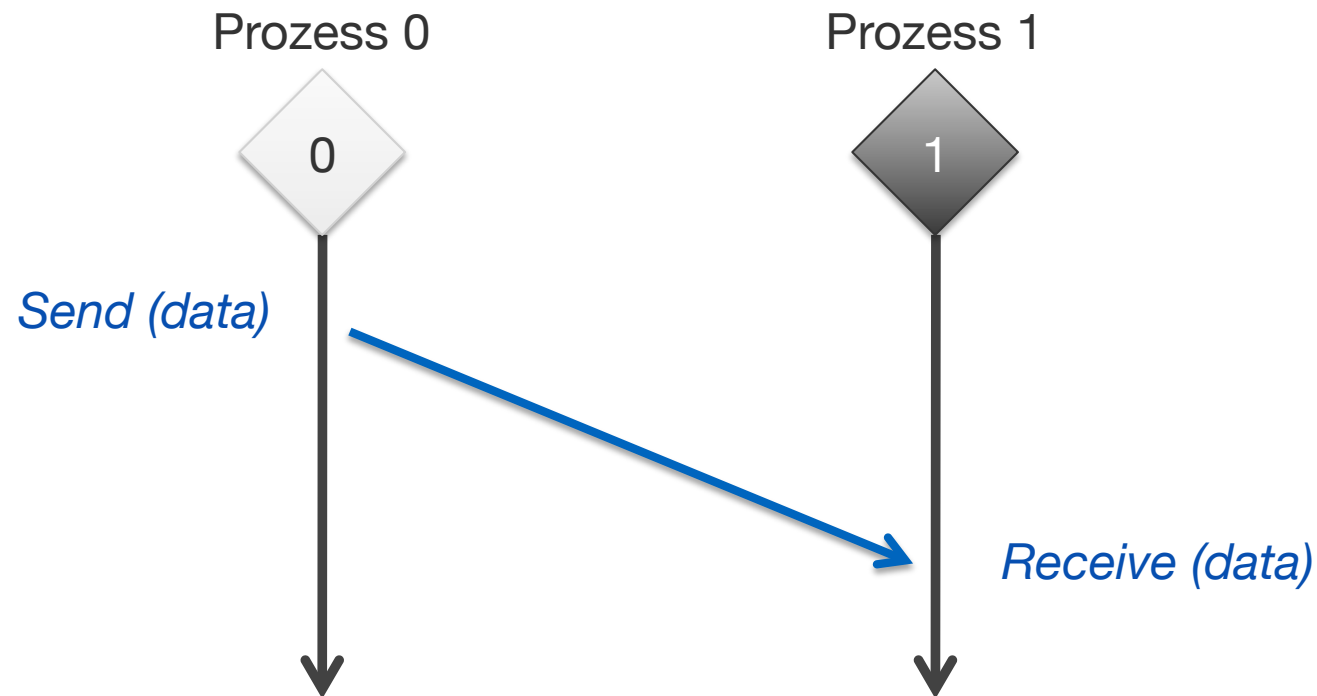
```
franz-diebold:Release Franz$ mpirun -np 5 MPI_HelloWorld  
Hallo Welt! Ich bin Prozess 0 von 5.  
Hallo Welt! Ich bin Prozess 3 von 5.  
Hallo Welt! Ich bin Prozess 1 von 5.  
Hallo Welt! Ich bin Prozess 2 von 5.  
Hallo Welt! Ich bin Prozess 4 von 5.
```

- Reihenfolge nicht deterministisch!

# Agenda

1. Einführung, Motivation
2. Kommunikationsmodell
3. *Punkt-Zu-Punkt*-Kommunikation
4. Globale Kommunikation
5. Vergleich MPI und OpenMP

### 3. *Punkt-Zu-Punkt-Kommunikation*



- Benötigt Kooperation beider Kommunikationspartner



### 3. *Punkt-Zu-Punkt*-Kommunikation

#### MPI-Funktionen zum Senden und Empfangen von Nachrichten:

- **MPI\_Send** (`void` \*message\_buffer, // zu verschickende Daten  
          `int` count, // „Länge“ der Daten  
          MPI\_Datatype type, // Datentyp der Daten  
          `int` destination, // Ziel-Prozess der Nachricht  
          `int` message\_tag, // Nachrichten-Identifikator  
          MPI\_Comm communicator); // Kontext
- **MPI\_Recv** (`void` \*message\_buffer, // zu empfangende Daten  
          `int` count, // „Länge“ der Daten  
          MPI\_Datatype type, // Datentyp der Daten  
          `int` source, // Herkunfts-Prozess der Nachricht  
          `int` message\_tag, // Nachrichten-Identifikator  
          MPI\_Comm communicator, // Kontext  
          MPI\_Status\* status); // Nachrichten-Status

## Datentypen in MPI (MPI\_Datatype)

MPI-Datentyp	C-Datentyp
MPI_SHORT	signed short int
MPI_INT	signed int
MPI_CHAR	signed char
MPI_LONG	signed long int
MPI_UNSIGNED_CHAR	unsigned char
MPI_UNSIGNED_SHORT	unsigned short
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long int
MPI_FLOAT	float
MPI_DOUBLE	double
MPI_LONG_DOUBLE	long double

### 3. *Punkt-Zu-Punkt*-Kommunikation

#### MPI-Funktionen zum Senden und Empfangen von Nachrichten:

- Nachricht von beliebigem Prozess:  
`MPI_ANY_SOURCE`
- Beliebige Nachricht empfangen:  
`MPI_ANY_TAG`
- Prüfen, ob zu empfangende Nachricht verfügbar:  
`MPI_Probe (int source, int message_tag,  
MPI_Comm comm, MPI_Status* status);`
  - Herkunft: `status->MPI_SOURCE`
  - Identifikator: `status->MPI_TAG`
  - Fehler: `status->MPI_ERROR`
  - „Länge“ der Daten: `MPI_Get_count (MPI_Status* status,  
MPI_Datatype type, int* count);`

### 3. *Punkt-Zu-Punkt*-Kommunikation

#### Beispiel: Producer-Consumer-MPI-Programm in C:

```
int numtasks, rank, producer, consumer, data, msgtag;
MPI_Status status;
producer = 0; consumer = 1; msgtag = 134;
[...]
```

```
if (rank == producer) {                                // Producer
    data = 1;
    while (data > 0) {
        scanf("%d", &data);
        MPI_Send(&data, 1, MPI_INT, consumer, msgtag, MPI_COMM_WORLD); }
} else if (rank == consumer) {                          // Consumer
    data = 1;
    while (data > 0) {
        MPI_Recv(&data, 1, MPI_INT, producer, msgtag, MPI_COMM_WORLD,
        &status); }
}
[...]
```

### 3. *Punkt-Zu-Punkt*-Kommunikation

#### Kombinierte Send- & Receive-Funktion:

- **MPI\_Sendrecv** (`void` \*send\_buffer, // zu verschickende Daten  
int count\_send // „Länge“ der Daten  
MPI\_Datatype send\_type, // Datentyp der Daten  
int destination, // Ziel-Prozess der Nachricht  
int send\_tag, // Nachrichten-Identifikator  
void \*rec\_buffer, // zu empfangende Daten  
int count\_rec, // „Länge“ der Daten  
MPI\_Datatype rec\_type // Datentyp der Daten  
int source, // Herkunfts-Prozess der Nachricht  
int rec\_tag, // Nachrichten-Identifikator  
MPI\_Comm communicator, // Kontext  
MPI\_Status\* status); // Nachrichten-Status

### 3. *Punkt-Zu-Punkt*-Kommunikation

#### Unterscheidung blockierend – nicht blockierend:

- **blockierend:**
  - Funktionsaufruf terminiert erst, wenn Puffer wieder verwendet werden bzw. gelesen werden darf
  - **MPI\_Send, MPI\_Recv, MPI\_Ssend, MPI\_Bsend, MPI\_Rsend, MPI\_Sendrecv**
- **nicht blockierend:**
  - Funktionsaufruf terminiert sofort (kein Warten auf Ende der Kommunikation)
  - **MPI\_Isend, MPI\_Issend, MPI\_Ibsend, MPI\_Irecv**
  - Puffer darf wieder verwendet werden: **MPI\_Wait, MPI\_Test, MPI\_Waitall**

### 3. *Punkt-Zu-Punkt*-Kommunikation

#### Nicht-blockierendes Senden & Empfangen:

- **MPI\_Isend** (`void` \*message\_buffer, // zu verschickende Daten  
int count, // „Länge“ der Daten  
MPI\_Datatype type, // Datentyp der Daten  
int destination, // Ziel-Prozess der Nachricht  
int message\_tag, // Nachrichten-Identifikator  
MPI\_Comm communicator // Kontext  
MPI\_Request\* request); // Request-Objekt (handle)
- **MPI\_Irecv** (`void` \*message\_buffer, // zu empfangende Daten  
int count, // „Länge“ der Daten  
MPI\_Datatype type, // Datentyp der Daten  
int source, // Herkunfts-Prozess der Nachricht  
int message\_tag, // Nachrichten-Identifikator  
MPI\_Comm communicator, // Kontext  
MPI\_Request\* request); // Request-Objekt (handle)

### 3. *Punkt-Zu-Punkt*-Kommunikation

#### Beispiel: Nicht-blockierende Kommunikation:

```
[...]
int rank, numtasks, left, right;
int buffer[10], buffer2[10];
MPI_Request request, request2;
MPI_Status status;
[...]
right = (rank+ 1) % numtasks;
left = rank- 1;
if (left < 0)
    left = numtasks - 1;
MPI_Irecv(buffer, 10, MPI_INT, left, 123, MPI_COMM_WORLD, &request);
MPI_Isend(buffer2, 10, MPI_INT, right, 123, MPI_COMM_WORLD, &request2);
MPI_Wait(&request, &status);
MPI_Wait(&request2, &status);
[...]
```

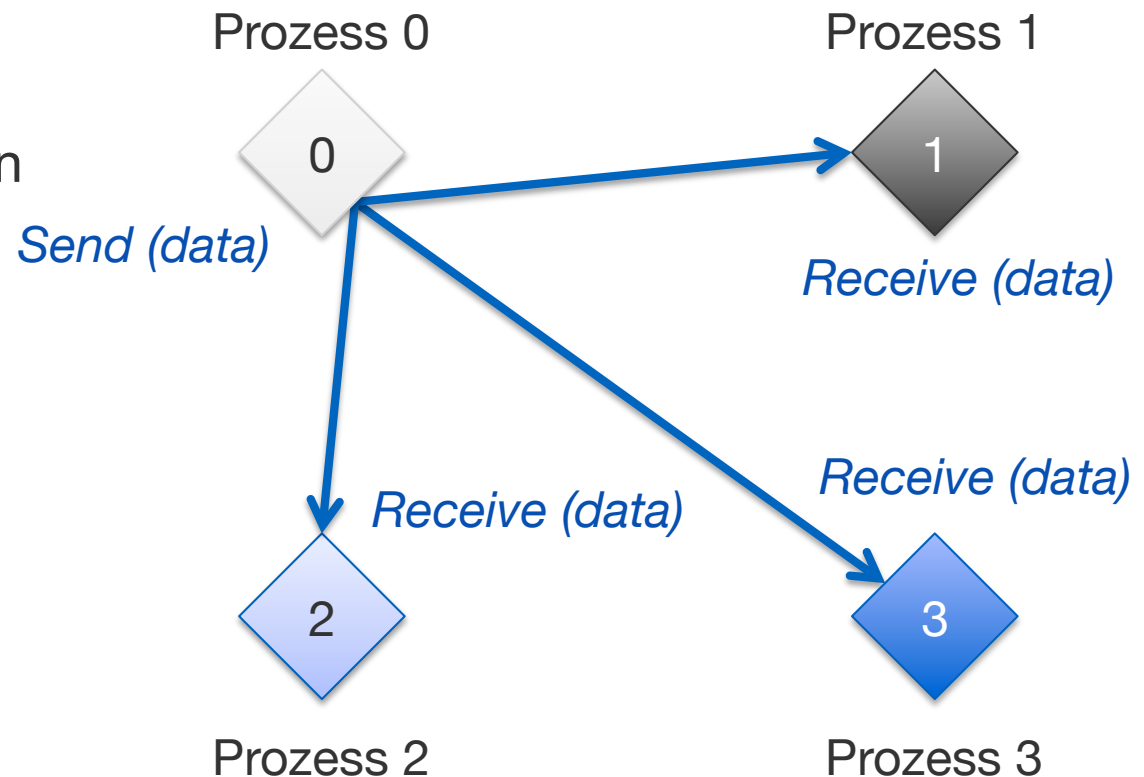


# Agenda

1. Einführung, Motivation
2. Kommunikationsmodell
3. *Punkt-Zu-Punkt*-Kommunikation
4. Globale Kommunikation
5. Vergleich MPI und OpenMP

## 4. Globale Kommunikation

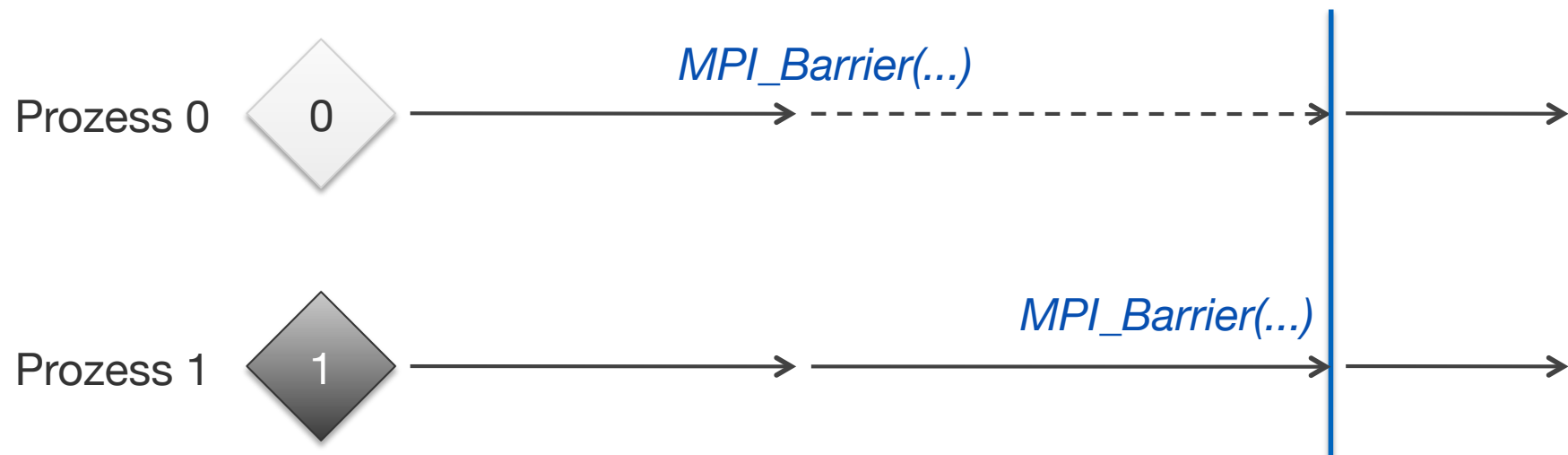
- Betrifft immer *alle* Prozesse des jeweiligen Kommunikators
- Typen:
  - Synchronisation
  - Datenübertragung
  - Gemeinsame Berechnungen
- Keine „Tags“ (Nachrichten-Identifikatoren)



## 4. Globale Kommunikation

### Barriere:

- Blockiert Prozesse, bis alle diese Funktion aufgerufen haben

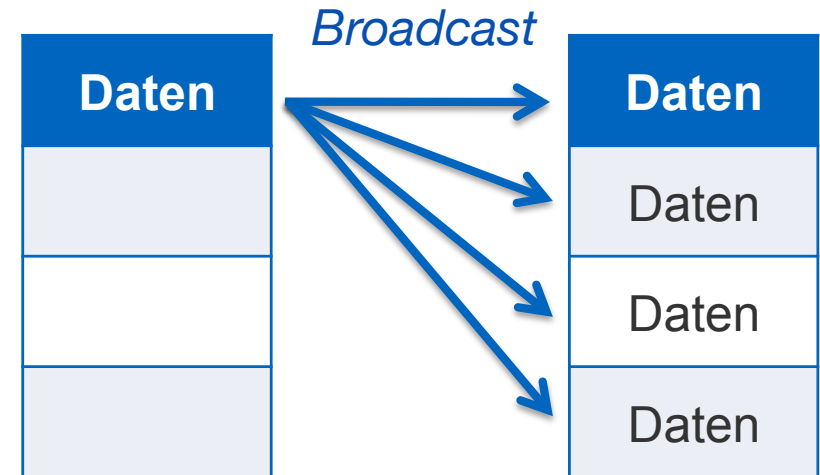


- **MPI\_Barrier** (MPI\_Comm comm); // Kommunikator (Kontext)

## 4. Globale Kommunikation

### Broadcast:

- Schickt Nachricht von „root“-Prozess zu allen anderen Prozessen
- Bei Sender und Empfänger gleicher Funktionsaufruf
- Barriere (vgl. **MPI\_Barrier**)

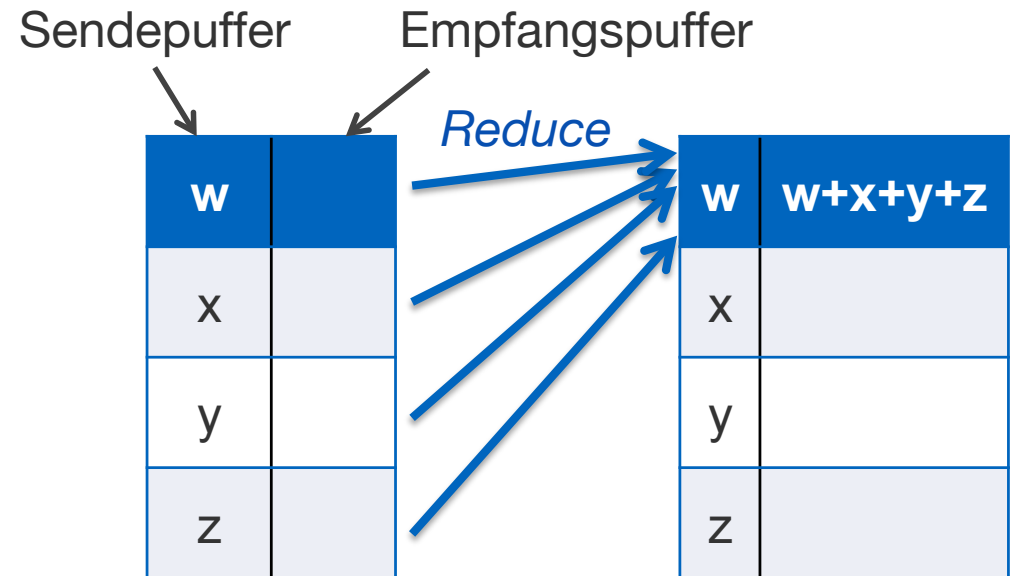


- **MPI\_Bcast** (`void*` buffer,  
                  `int` count,  
                  MPI\_Datatype type,  
                  `int` root,  
                  MPI\_Comm comm);  
                  // Sende-/Empfangspuffer  
                  // Nachrichtenlänge  
                  // Datentyp  
                  // „root“-Prozess  
                  // Kommunikator (Kontext)

## 4. Globale Kommunikation

### Reduce:

- Führt globale Operation aus
- „root“-Prozess erhält Ergebnis
- Bei Sender und Empfänger gleicher Funktionsaufruf



- **MPI\_Reduce** (`void*` sendbuffer, `void*` recvbuffer, `int` count, `MPI_Datatype` type, `MPI_Op` op, `int` root, `MPI_Comm` comm);
  - // Sendepuffer
  - // Empfangspuffer
  - // Nachrichtenlänge
  - // Datentyp
  - // Reduktionsfunktion
  - // „root“-Prozess
  - // Kommunikator

## Globale Operationen in MPI (MPI\_Op)

MPI-Op	Bedeutung
MPI_MAX	Maximum
MPI_MIN	Minimum
MPI_SUM	Summe
MPI_PROD	Produkt
MPI_LAND / MPI_LOR / MPI_LXOR	logisches AND / OR / XOR
MPI_BAND / MPI_BOR / MPI_BXOR	bitweises AND / OR / XOR
MPI_MAXLOC	Maximum und Index des Prozesses
MPI_MINLOC	Minimum und Index des Prozesses

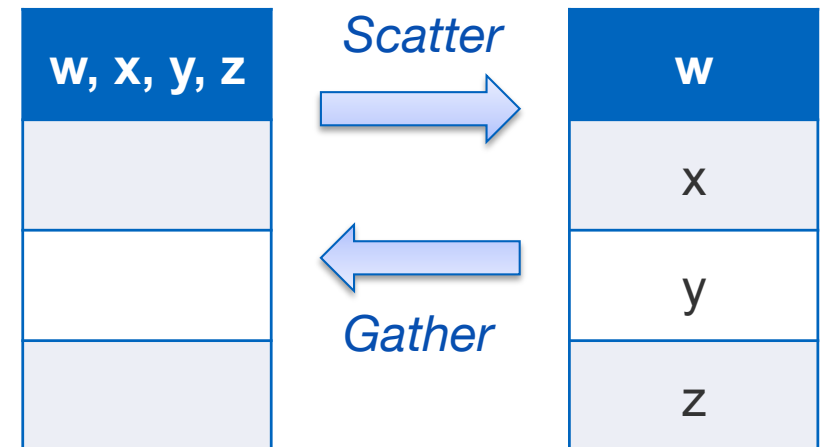
- Nutzung eigener Reduktionsfunktionen möglich mit:

```
MPI_Op_create (MPI_User_function *function, int commute, MPI_Op *op);
```

## 4. Globale Kommunikation

### Gather / Scatter:

- Daten „einsammeln“ bzw. „verteilen“
- Bei Sender und Empfänger gleicher Funktionsaufruf



- **MPI\_Gather/** `(void* sendbuffer,` // Sendepuffer  
**MPI\_Scatter** `int send_count,` // Nachrichtenlänge  
`MPI_Datatype send_type,` // Datentyp (Senden)  
`void* recvbuffer,` // Empfangspuffer  
`int send_count,` // Nachrichtenlänge  
`MPI_Datatype send_type,` // Datentyp (Empfang)  
`int root,` // „root“-Prozess  
`MPI_Comm comm);` // Kommunikator

# Agenda

1. Einführung, Motivation
2. Kommunikationsmodell
3. *Punkt-Zu-Punkt*-Kommunikation
4. Globale Kommunikation
5. Vergleich MPI und OpenMP



## 5. Vergleich MPI und OpenMP

MPI	OpenMP
✓ mit verteiltem und gemeinsamem Speicher möglich	✓ einfacher zu programmieren und debuggen als MPI
✓ mehr Möglichkeiten als OpenMP	✓ Programm auch seriell ausführbar
✓ Rechner mit verteiltem Speicher billiger als mit gemeinsamem Speicher	✓ Direktiven inkrementell hinzufügbare
	✓ Code einfacher zu verstehen und zu warten
✗ Flaschenhals: Netzwerk	✗ nur mit gemeinsamem Speicher nutzbar
✗ schwer zu debuggen	✗ Compiler mit OpenMP-Unterstützung notwendig
✗ mehr Änderungen von seriell zu parallel	(✗ hauptsächlich für Schleifenparallelisierung)

## Quellen

- <http://www.tu-chemnitz.de/informatik/RA/projects/mpihelp/>
- <https://computing.llnl.gov/tutorials/mpi/>
- <http://www.mpi-forum.org/docs/mpi21-report.pdf>
- [http://www.rz.rwth-aachen.de/global/show\\_document.asp?id=aaaaaaaaaablubl](http://www.rz.rwth-aachen.de/global/show_document.asp?id=aaaaaaaaaablubl)