

## 1 MPI: Message Passing Interface

- Portable Bibliothek zum Nachrichtenaustausch
- Wurde 1993-94 durch ein intern. Gremium entwickelt
- 1997 wurde der Standard überarbeitet (MPI2)
- Open-Source Implementierungen gibt beispielweise von: **MPICH**<sup>1</sup> und **OpenMPI**<sup>2</sup>
- Eigenschaften von MPI:
  - Direkte Anbindung an C, C++ und Fortran
  - verschiedene Arten von Punkt-zu-Punkt Kommunikation
  - globale Kommunikation
  - Daten Umwandlung in heterogenen Systemen
  - Vier Netzwerke & Topologien möglich

### 1.1 Beispiel

#### Hello World

```
1 #include <mpi.h>
2 #include <iostream>
3
4 int main(int argc, char **argv)
5 {
6     int *buf, i, rank, nints, len;
7     char hostname[256];
8
9     MPI_Init(&argc,&argv);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11    gethostname(hostname,255);
12    std::cout << "Hello world! I am process number:" << rank
13    << " on host " << hostname << std::endl;
14    MPI_Finalize();
15    return 0;
16 }
```

#### Kompilieren eines Programms

- Das Beispiel ist im SPMD-Stil geschrieben; der MPI Standard legt dies nicht fest und man kann auch andere Paradigmen verwenden.
- Kompilieren eines MPI C-Programmes und starten mit 8 Prozessen:

```
mpicc -o hello hello.c
mpirun -machinefile machines -np 8 hello
```

Die Liste der Computer steht in der Datei machines.

C für C++ Programme heißt der Compiler

```
mpicxx -o hello hello.cc
```

---

<sup>1</sup><http://www-unix.mcs.anl.gov/mpi/mpich>

<sup>2</sup><http://www.open-mpi.org/>

## Initialisierung und Beenden

Bevor irgend ein MPI Befehl aufgerufen werden darf, muss MPI mit `MPI_Init` initialisiert werden, damit MPI das parallele Programm starten kann.

```
int MPI_Init(int *argc, char ***argv)
```

Nach dem letzten MPI Aufruf wird `MPI_Finalize` ausgeführt, um alles Prozesse ordentlich zu beenden.

```
int MPI_Finalize(void)
```

## 1.2 Kommunikatoren und Topologien

MPI erlaubt es Kommunikation auf einem Subset der gestarteten Prozesse durchzuführen, indem *virtuelle* Netzwerke, sogenannte Kommuikatoren, angelegt werden.

- `MPI_Comm` beschreibt einen Kommunikator, eine Menge von Prozessen  $\{0, \dots, P - 1\}$ .
- Das vordefinierte Netzwerk `MPI_COMM_WORLD` enthält alle gestarteten Prozesse.
- *Virtuelle Topologien*: Ein Kommunikator kann zusätzlich eine spezielle Struktur erhalten, z.B. ein mehrdimensionales Feld, oder ein allgemeiner Graph.
- *Kontext*: Jeder Kommunikator definiert seinen eigenen Kommunikationskontext.

## Rank und Size

Die Anzahl der Prozesse in einem Kommunikator wird mit `MPI_Comm_size` bestimmt:

```
int MPI_Comm_size(MPI_Comm comm, int *size)
```

Innerhalb eines Kommunikators hat jeder Prozess eine eindeutige Nummer, diese wird mit `MPI_Comm_rank` bestimmt:

```
int MPI_Comm_rank(MPI_Comm comm, int *rank)
```

```
1 #include <stdio.h>
2 #include <string.h>
3 #include <mpi.h> // provides MPI macros and functions
4
5     int main (int argc, char *argv[])
6 {
7     int my_rank;
8     int P;
9     int dest;
10    int source;
11    int tag=50;
12    char hostname[256];
13    MPI_Status status;
14
15    MPI_Init(&argc,&argv); // begin of every MPI program
16
17    MPI_Comm_size(MPI_COMM_WORLD ,&P);           // number of processes
18    MPI_Comm_rank(MPI_COMM_WORLD ,&my_rank); // my process number
19
20    gethostname(hostname ,255);
21
22    // number of current process always between 0 and P-1
23    if (my_rank!=0)
24    {
25        dest = 0;
26        MPI_Send(hostname ,strlen(hostname)+1,MPI_CHAR , // Send data
27                  dest ,tag ,MPI_COMM_WORLD); // (blocking)
28    }
29 }
```

```

    else
31   {
32     for (source=1; source<P; source++)
33     {
34       MPI_Recv(hostname ,256,MPI_CHAR,source ,tag , // Receive data
35           MPI_COMM_WORLD ,&status); // (blocking)
36       std::cout << "Received a message from process "
37           << "on machine " << hostname;
38     }
39   }
40
41   MPI_Finalize(); // end of every MPI program
42
43   return 0;
}

```

### Ausgabe des Beispiel Programms (mit P=8)

```

I am process 0 of 8
2
I am process 1 of 8
4
I am process 2 of 8
6
I am process 3 of 8
8
I am process 4 of 8
10
I am process 5 of 8
12
I am process 6 of 8
14
I am process 7 of 8

```

### 1.3 Blockierende Kommunikation

Die Entsprechung zu **send** und **recv** bieten

```

int MPI_Send(void *message, int count, MPI_Datatype dt,
             int dest, int tag, MPI_Comm comm);
int MPI_Recv(void *message, int count, MPI_Datatype dt,
             int src, int tag, MPI_Comm comm,
             MPI_Status *status);

```

Die ersten drei Parameter, **message**, **count** und **dt**, beschreiben die eigentlichen Daten. **message** ist ein Zeiger auf ein Feld mit **count** Elementen des Typs **dt**. Die Angabe des Datentyps erlaubt die automatische Umwandlung durch MPI. Die Parameter **dest**, **tag** und **comm** beschreiben das Ziel, bzw. die Quelle der Nachricht.

MPI biete verschiedene Varianten von **MPI\_Send** (**MPI\_BSend**, **MPI\_SSend**, **MPI\_RSend**), die wir aber jetzt nicht weiter diskutieren wollen.

**MPI\_ANY\_SOURCE** und **MPI\_ANY\_TAG** können verwendet werden, um beliebige Nachrichten zu empfangen. Damit enthält **MPI\_Recv** die Funktionalität von **recv\_any**.

## Datenumwandlung

MPI erlaubt die Verwendung in heterogenen Netzen. Hierbei ist es nötig manche Daten an die Darstellung auf der fremden Architektur anzupassen.

MPI definiert die architektur-unabhängigen Datentypen:

MPI\_CHAR, MPI\_UNSIGNED\_CHAR, MPI\_BYTE MPI\_SHORT, MPI\_INT, MPI\_LONG, MPI\_LONG\_LONG\_INT, MPI\_UNSIGNED, MPI\_UNSIGNED\_SHORT, MPI\_UNSIGNED\_LONG, MPI\_FLOAT, MPI\_DOUBLE and MPI\_LONG\_DOUBLE.

Der MPI Datentyp MPI\_BYTE wird *nie* konvertiert.

## Status

```
1 typedef struct {
2     int count;
3     int MPI_SOURCE;
4     int MPI_TAG;
5     int MPI_ERROR;
} MPI_Status;
```

MPI\_Status ist ein zusammengesetzter Datentyp, der Informationen über die Anzahl der empfangenen Objekte, den Quellprozess, das Tag und den Fehler Status enthält.

## Guard Funktion

Die Guard Funktion **rprobe** liefert

```
1 int MPI_Iprobe(int source, int tag, MPI_Comm comm,
2                 int *flag, MPI_Status *status);
```

Es ist eine nicht-blockierende Funktion, die überprüft, ob eine Nachricht vorliegt. flag erhält den Wert **true** ( $\neq 0$ ) wenn eine Nachricht mit passendem source und tag empfangen werden kann. Auch hier können MPI\_ANY\_SOURCE und MPI\_ANY\_TAG verwendet werden.

## 1.4 Nicht-blockierende Kommunikation

Die Funktionen **asend** und **arecv** bietet MPI als

```
1 int MPI_ISend(void *buf, int count, MPI_Datatype dt,
2                int dest, int tag, MPI_Comm comm,
3                MPI_Request *req);
4 int MPI_IRecv(void *buf, int count, MPI_Datatype dt,
5               int src, int tag, MPI_Comm comm,
6               MPI_Request *req);
```

MPI\_Request speichert den Status einer Kommunikation, wie unsere **msgid**.

### MPI\_Request-Objekte

Der Status einer Nachricht kann mit Hilfe der MPI\_Request-Objekte und folgender Funktion geprüft werden:

```
int MPI_Test(MPI_Request *req, int *flag, MPI_Status *status);
```

flag wird auf **true** ( $\neq 0$ ) gesetzt, wenn die Kommunikation, die req beschreibt abgeschlossen ist. In diesem Fall enthält status weiter Informationen.

## 1.5 Globale Kommunikation

MPI bietet ebenfalls Funktionen zur globalen Kommunikation, welche alle Prozesse eines Kommunikators einschließen.

```
1 int MPI_Barrier(MPI_Comm comm);
```

implementiert eine Barriere; alle Prozesse werden blockiert, bis der letzte Prozess die Funktion ausgeführt hat.

```
1 int MPI_Bcast(void *buf, int count, MPI_Datatype dt,
               int root, MPI_Comm comm);
```

verteilt eine Nachricht an alle Prozesse eines Kommunikators (Einer-an-Alle Kommunikation).

## Einsammeln von Daten

MPI hat eine Reihe verschiedener Funktionen um Daten von verschiedenen Prozessen einzusammeln. Z.B:

```
1 int MPI_Reduce(void *sbuf, void *rbuf, int count,
                MPI_Datatype dt, MPI_Op op, int root, MPI_Comm comm);
```

kombiniert die Daten im Sende-Puffer `sbuf` aller Prozesse durch die assoziative Operation `op` (z.B. `MPI_SUM`, `MPI_MAX` oder `MPI_MIN`). Das Ergebnis erhält der Prozesse `root` in seinen Empfang-Puffer `rbuf`.

## 2 Zeitmessung

MPI bietet direkt Funktionen zur Zeitmessung:

```
double MPI_Wtime();
```

Der Rückgabewert sind Sekunden seit einem „beliebigen“ Zeitpunkt in der Vergangenheit. Die Zeit für eine spezielle Operation lässt sich also mit

```
1 double start = MPI_Wtime();
2 expensive_funktion();
3 std::cout << "elapsed time = " << MPI_Wtime() - start << std::endl;
```

bestimmen.

## References

- [1] *MPI: Dokumentation der verschiedenen Message-Passing Interface Standards* <http://www mpi-forum.org/docs/>
- [2] *MPICH-A Portable Implementation of MPI* <http://www-unix.mcs.anl.gov/mpi/mpich>
- [3] *Open MPI: Open Source High Performance Computing* <http://www.open-mpi.org/>
- [4] *Liste von MPI Tutorials* <http://www-unix.mcs.anl.gov/mpi/tutorial/>