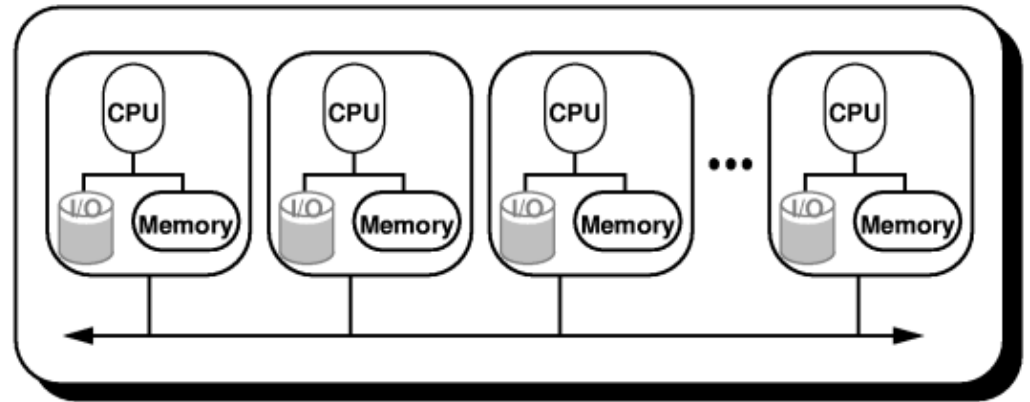


**MPI**  
**Message Passing**  
**Interface**

- **Autonómne procesory s vlastnou pamäťou** prepojené komunikačnou sieťou



- Komunikácia realizovaná posielaním správ
- Procesory sú typicky od seba fyzicky vzdialené (clustre, BOINC, „superpočítače“, ...)

- **API**, ktoré umožňuje procesom komunikovať prostredníctvom **odosielania a prijímania správ**
- používa sa pri vývoji paralelných programov pre superpočítače a clustre
- ciele:
  - použiteľnosť, portabilita, efektívnosť, flexibilita
- vznikla v roku **1994** ako výsledok dohody 40 organizácií združených v MPI fórum
- <http://www.mcs.anl.gov/research/projects/mpi/>

- **MPI nie je implementácia**, resp. nejaká knižnica, je to len špecifikácia API
- štandardizované API umožňuje vývoj portabilných paralelných programov
- výrobcovia hardvéru vytvárajú implementácie MPI **optimalizované** pre svoj hardvér
- známe implementácie:
  - MPICH, LAM/MPI, WMPI, OpenMPI, Intel MPI
  - MS MPI (optimalizované pre Win HPC Server 2008)

- objektovo-orientované Java rozhranie pre štandardizované MPI (1997)
- <http://www.hpjava.org/mpiJava.html>
- implementácie:
  - **mpiJava** - java wrapper (cez JNI) na MPICH implementáciu
  - **MPJ Express** - čistá Java implementácia
    - <http://mpj-express.org/>
  - **F-MPJ**

- implementácia java MPI rozhrania len v čistej Jave
  - <http://mpj-express.org/>
- 2 konfigurácie:
  - multicore konfigurácia
  - cluster konfigurácia
    - cez Java NIO
    - cez Myronet
- ukážka inštalácie...

```
import mpi.*;

public class HelloWorld {
    public static void main(String[] args) {
        MPI.Init(args);
        int size = MPI.COMM_WORLD.Size();
        int me = MPI.COMM_WORLD.Rank();
        System.out.println("Hi from <"+me+">");
        MPI.Finalize();
    }
}
```

- počet procesov je určený pri spustení programu (argument -np)

- **Skupina** (group)
  - nejaká podmnožina vytvorených procesov
- **Komunikátor** (Comm)
  - umožňuje komunikáciu v rámci skupiny procesov
  - každý proces je v rámci komunikátora identifikovaný svojim **rank**-om
  - **MPI.COMM\_WORLD** je komunikátor pre skupinu tvorenú všetkými procesmi



```
char[] buffer = ...
```

```
MPI.COMM_WORLD.Send(  
  buffer,  
  offset,  
  count,  
  MPI.CHAR,  
  target,  
  tag);
```

Pole hodnôt  
primitívneho typu

Offset, od ktorého  
začínajú dátové elementy,  
ktoré sa majú odoslať

Počet odosielaných  
dátových elementov

Typ odosielaných dátových  
elementov

Rank adresáta v  
komunikátore

Tag správy

| MPI datatype | Java datatype |
|--------------|---------------|
| MPI.BYTE     | byte          |
| MPI.CHAR     | char          |
| MPI.SHORT    | short         |
| MPI.BOOLEAN  | boolean       |
| MPI.INT      | int           |
| MPI.LONG     | long          |
| MPI.FLOAT    | float         |
| MPI.DOUBLE   | double        |
| MPI.PACKED   |               |

*Status Comm.Recv(Object buf,*

*int offset,*

*int count,*

*Datatype datatype,*

*int source,*

*int tag)*

Prijímací buffer je určený ako časť (určená offsetom a parametrom count) poľa hodnôt primitívneho typu

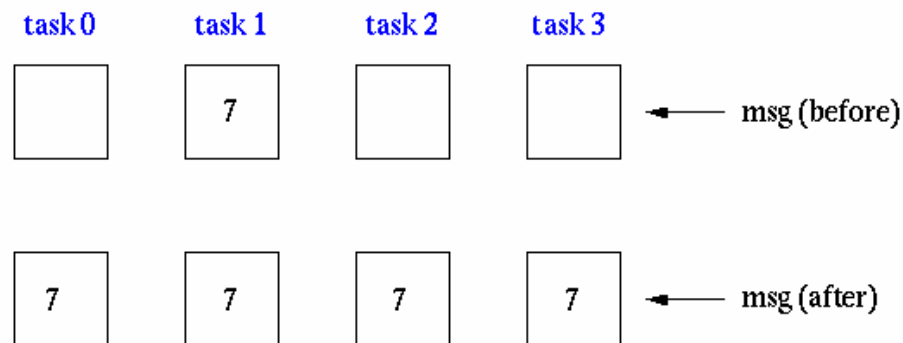
Rank odosielateľa (alebo MPI.SOURCE\_ANY)

Tag prijímanej správy (alebo MPI.TAG\_ANY)

Zo Status objektu sa vieme dozvedieť skutočný počet prijatých dátových elementov, odosielateľa a tag správy.

- **blokované vs. neblokované**
  - pri blokovaných operáciách sa čaká na dokončenie operácie (napr. pri odosielaní sa čaká, kým správa bude odoslaná; po ukončení blokovanej operácie je možné znovupoužívať buffre) ... (napr. MPI\_Isend)
- **synchrónne vs. asynchrónne**
  - pri synchrónnych operáciách máme potvrdenie prijatia správy príjemcom (napr. MPI\_Ssend)
- Send a Recv sú asynchrónne blokované.

- **MPI\_Barrier** (Barrier)
  - blokuje vykonávanie procesu dokiaľ všetky procesy v komunikátore nevykonajú túto operáciu
- **MPI\_Bcast** (Bcast)
  - pošle správu všetkým procesom v komunikátore, resp. prijme broadcastovanú správu



- Scatter
- Gather
- Allgather
- Reduce
- Allreduce
- Alltoall
- Scan

- MPI umožňuje vytvárať **vlastné skupiny** z existujúcich skupín
  - ku každej skupine možno vytvoriť komunikátor, ktorý je užitočný v prípade kolektívnej komunikácie
- MPI umožňuje definovať **virtuálne topológie** procesov
  - mriežka, graf
  - MPI implementácia môže túto informáciu využiť na mapovanie procesov na fyzické procesory

Ďakujem za pozornosť !

