

## MPI (Message Passing Interface)

- rozhraní pro (explicitní) programování paralelních aplikací na jednom i více počítačích s **distribuovanou pamětí**
- startuje **několik procesů** (každý s nezávislou pamětí) zasílajících si za chodu zprávy (data)
- zdrojový kód procesů obvykle totožný (**SPMD**, single program multiple data), možné i schéma master-slave (**MPMD**)
- realizace: podprogramy, funkce, proměnné a konstanty, proměnné prostředí
- překlad pomocí wrapperů (př. **mpif90**) zajišťujících připojení knihoven, spouštění pomocí agenta (**mpirun/mpiexec**)
- <http://www.mpi-forum.org>, dokumentace **MPI: A Message-Passing Interface Standard ver. 2.2** (PDF 647 stran)
- 1994 standard **MPI-1** (2008: ver. 1.3), 1996 **MPI-2** (2009: ver. 2.2)
- implementace: knihovny (včetně zdrojových kódů) **Open MPI** (ver. 1.4 s MPI-2), **MPICH** aj. pro **Fortran**, **C/C++** aj.

### Součásti MPI

- MPI-1 – dvouprocesová komunikace (**point-to-point communication**: send, receive)
- kolektivní komunikace (**collective communication**: broadcast, scatter, gather, reduce)
  - **virtuální topologie** procesů (kartézská, grafová)
  - **synchronizace** procesů (blokující a neblokující komunikace, bariéry)
  - **informace** o procesech a jejich skupinách (comm, rank, size)
- MPI-2 – jednostranná komunikace pro přístup do vzdálené paměti (one-sided communication: put, get, accumulate)
- dynamické řízení procesů pro schéma master-slave (dynamic process management)
  - ošetření vstupu a výstupu na distribuovaných systémech (MPI-IO)

### Použití MPI ve Fortranu

Pojmy: komunikátor = skupina komunikujících procesů (**communicator**), identifikace procesu (**rank**), „procesor“ = proces

Časté argumenty procedur (buf standardního nebo definovaného typu, ostatní integer):

**ierror**: chybový kód

**comm**: číslo komunikátoru, komunikátor **MPI\_COMM\_WORLD** sdružuje všechny procesy

**size**: velikost komunikátoru

**rank**: číslo procesu

**buf, sendbuf, recvbuf**: proměnná (skalár, pole, struktura) standardního nebo definovaného typu

**count**: počet prvků v poli buf, 1 pro skalár

**datatype**: typ proměnné buf, př. konstanty **MPI\_INTEGER**, **MPI\_REAL**, **MPI\_REAL8**, **MPI\_COMPLEX** aj., **MPI\_CHARACTER** odpovídá typu CHARACTER(1), víceznačkové řetězce jsou ekvivalentní poli znaků, MPI neprovádí typové konverze (př. INTEGER na REAL), provádí však reprezentační konverze (př. default REAL na default REAL)

**source**: číslo zdrojového procesu v komunikaci dvojice, může být **MPI\_ANY\_SOURCE**

**dest**: číslo cílového procesu v komunikaci dvojice

**root**: číslo zdrojového nebo cílového procesu v kolektivní komunikaci

**tag**: číslo zprávy, může být **MPI\_ANY\_TAG**

**status**: integer pole o velikosti **MPI\_STATUS\_SIZE**, může být **MPI\_STATUS\_IGNORE**

**request**: číslo požadavku při neblokující komunikaci

MPI procedury ve Fortranu jsou implementovány převážně jako podprogramy, existují i funkce a pojmenované konstanty. MPI v C má podobu funkcí vracejících ierror (oproti Fortranu nemají argument ierror).

### Inicializace a finalizace, informace

Deklarace MPI jmen zajišťuje vložení souboru mpif.h. Před použitím je třeba knihovnu MPI inicializovat a je vhodné zjistit číslo a velikost komunikátoru a číslo procesu; po použití lze knihovnu finalizovat.

- deklarace `include 'mpif.h'` ! může připojovat další soubory: mpi-config.h, mpi-common.h aj.
- inicializace `call MPI_INIT(ierr)`
- výchozí komunikátor `comm=MPI_COMM_WORLD`
- velikost komunikátoru `call MPI_COMM_SIZE(comm,size,ierr)`
- číslo procesu `call MPI_COMM_RANK(comm,rank,ierr)`
- finalizace `call MPI_FINALIZE(ierr)`

## Dvouprocesová komunikace

Odesílatel odesílá jednu zprávu jednomu příjemci ve stejném komunikátoru. Zprávy se skládají z obálky (source, dest, tag, comm) a dat (proměnná buf typu datatype, je-li polem, pak o count prvcích). Přijímaná zpráva nesmí být delší než očekávaná, může být kratší (její aktuální délku vrací MPI\_GET\_COUNT). Příjemce může přijmout zprávu od konkrétního nebo libovolného odesílatele. Zprávy mezi tímž odesílatelem a příjemcem se nepředbíhají.

Odeslání zprávy (send): **MPI\_SEND**(buf,count,datatype,dest,tag,comm,ierror)

Příjem zprávy (receive): **MPI\_RECV**(buf,count,datatype,source,tag,comm,status,ierror)

Skutečná délka zprávy: **MPI\_GET\_COUNT**(status,datatype,count,ierror)

Nezáleží-li příjemci na konkrétním source, tag a status, lze užít předdefinovaných MPI\_ANY\_SOURCE, MPI\_ANY\_TAG a MPI\_STATUS\_IGNORE. Komunikace neproběhne, užije-li se prázdný dest nebo source, tj. MPI\_PROC\_NULL.

## Komunikační módy

**Blokující komunikace:** příjemce vždy čeká v proceduře pro příjem (receive) až do přijetí celé zprávy, odesílatel vždy čeká v proceduře pro odeslání (send) na odeslání celé zprávy, což není moment totožný s momentem přijetí zprávy příjemcem. Send je obvykle dokončen dříve než párový receive, není však vyloučen opak. Může být úspornější co do času i paměti, je-li receive spuštěn dříve než send.

Módy odesílání:

- **bufEROVANÝ mód (B):** odesílatel zkopíruje zprávu ze své paměti do send-bufu a může pokračovat v běhu nezávisle na připravenosti příjemce, implicitní velikost send-bufu (př. 32 KB) lze explicitně zvětšit, při nedostatečné velikosti send-bufu chyba (buffer overflow),
- **SYNCHRONNÍ mód (S):** odesílatel vyčká na připravenost příjemce (ušetří se bufery a kopírování mezi nimi) a pokračuje až po odeslání zprávy, nebude-li příjemce připraven nikdy, proces uvázne (deadlock),
- **READY mód (R):** odesílatel předpokládá připravenost příjemce (ušetří se inicializační fáze přenosu), při nepřipraveném příjemci nedefinovaný stav (např. čekání jako v synchronním módu),
- **STANDARDNÍ mód:** MPI volí podle situace (tj. v závislosti na připravenosti příjemce a dostupnosti buferů), kdy užije buferovaný a kdy synchronní mód.

Příslušné sendy (**MPI\_BSEND**, **MPI\_SSEND**, **MPI\_RSEND**) mají rozhraní totožné s MPI\_SEND a párují se s MPI\_RECV, který blokuje příjemce až do přijetí celé zprávy. K přenositelnosti programu přispívá, lze-li každý standardní MPI\_SEND nahradit synchronním MPI\_SSEND, tj. je-li program nezávislý na dostupnosti buferů v daném systému. Při ověřování dostupnosti buferů může být užitečné dočasně nahradit standardní MPI\_SEND buferovanými MPI\_BSEND. Odesílatel může připojit prostor pro send-bufer explicitně pomocí **MPI\_BUFFER\_ATTACH**(buffer,size,ierror) a později odpojit pomocí **MPI\_BUFFER\_DETACH**; to nemusí být relevantní pro standardní mód.

**Neblokující komunikace (I):** send i receive mohou být rozděleny na volání zahajovacího a kompletačního podprogramu (send-start a send-complete, receive-start a receive-complete), mezi kterými může proces vykonávat jinou práci. Mezi voláním send-start a send-complete by odesílatel neměl přepsat odesílanou proměnnou, mezi voláním receive-start a receive-complete by příjemce neměl přijímanou proměnnou ani číst, ani do ní psát. Ukončené odeslání neznamená, že byl dokončen příjem, ukončený příjem neznamená, že bylo dokončeno odeslání. Neblokující volání lze párovat s blokujícími voláními. Neblokující send-start lze kombinovat s módy B, S i R.

**Vytvoření požadavku** na odeslání nebo příjem:

– odeslání (send-start): **MPI\_ISEND**(buf,count,datatype,dest,tag,comm,request,ierror)

– příjem (receive-start): **MPI\_IRecv**(buf,count,datatype,source,tag,comm,request,ierror)

kde request je integer číslo požadavku (ukazatel na objekt) k užití následnými dotazy a při send-complete, varianty: MPI\_IBSEND, MPI\_ISSEND, MPI\_IRSEND.

**Kompletace požadavku** na odeslání nebo příjem:

– odeslání i příjem (send- i receive-complete):

blokující **MPI\_WAIT**(request,status,ierror)

neblokující **MPI\_TEST**(request,flag,status,ierror)

v případě ukončené operace provedou příslušné dealokace a deasociaci requestu a nastaví flag=.true.,

jinak MPI\_WAIT čeká na dokončení operace, zatímco MPI\_TEST nastaví flag=.false. a pokračuje,

neblokující **MPI\_REQUEST\_FREE**(request,ierror)

provede dealokaci a deasociaci requestu a ponechá možnost (netestovatelného) dokončení operace na pozadí,

varianty pro synchronizaci více požadavků (array\_of\_requests):

MPI\_WAITANY, MPI\_TESTANY, MPI\_WAITALL, MPI\_TESTALL, MPI\_WAIT SOME, MPI\_TESTSOME.

Ostatní podprogramy a funkce:

– zjištění informací o požadavku: MPI\_REQUEST\_GET\_STATUS, MPI\_PROBE, MPI\_IProbe,

– zrušení požadavku: MPI\_CANCEL.

**Perzistentní komunikace** pro opakující se přenosy s týmž seznamem argumentů:

**MPI\_SEND\_INIT** (s variantami B, S a R) a **MPI\_RECV\_INIT**; přenosy zahajuje volání **MPI\_START**(request), kompletuje **MPI\_WAIT** a **MPI\_TEST** (bez dealokace requestu), požadavek dealokuje **MPI\_REQUEST\_FREE**.

**Kombinace MPI\_SENDRECV** pro odeslání a příjem dvou zpráv mezi dvěma procesy, případně **MPI\_SENDRECV\_REPLACE** pro vzájemnou výměnu jedné zprávy.

### Kolektivní komunikace

Určeno pro hromadnou komunikaci uvnitř skupiny nebo mezi skupinami. Vedle **synchronizační procedury** (**BARRIER**) sem patří komunikace typu **one-to-all** (**BCAST**, **SCATTER**), **all-to-one** (**GATHER**, **REDUCE**) a **all-to-all** (**ALLTOALL** aj.). Nahrazuje sekvence dvouprocesové komunikace.

- synchronizace: **MPI\_BARRIER**(comm)
- rozeslání zprávy: **MPI\_BCAST**(buf,count,datatype, root,comm,ierror)
- rozeslání více zpráv: **MPI\_SCATTER**(sendbuf,sendcount,sendtype, recvbuf,recvcount,recvtype, root,comm,ierror)
- shromáždění zpráv: **MPI\_GATHER**(sendbuf,sendcount,sendtype, recvbuf,recvcount,recvtype, root,comm,ierror)
- redukce: **MPI\_REDUCE**(sendbuf,recvbuf,count,datatype, op, root,comm,ierror)  
kde op jsou konstanty **MPI\_SUM**, **MPI\_PROD**, **MPI\_LAND**, **MPI\_LOR**, **MPI\_MIN**, **MPI\_MINLOC** aj.

### Definice obecných datových typů

- pole **MPI\_TYPE\_CONTIGUOUS**(count,oldtype,newtype)
- struktura **MPI\_TYPE\_CREATE\_STRUCT**  
(count,array\_of\_blocklengths,array\_of\_displacements,array\_of\_types, newtype)
- pro (nepovinné) sbalení nespojitých dat do spojitěho buferu  
**MPI\_PACK**(inbuf,incount,datatype,outbuf,outsized,position,comm)  
**MPI\_UNPACK**(inbuf,insize,position,outbuf,outcount,datatype,comm)

### Virtuální topologie procesů

- intrakomunikátory: pro komunikaci uvnitř skupiny procesů
- interkomunikátory: komunikace mezi skupinami

### Překladový wrapper mpif90

Volá zvolený fortranský překladač a připojuje potřebné knihovny.

- překlad: **mpif90** source.f90
- informace o použitém překladači a nastavení: **mpif90 --showme**, **ompi\_info**
- nastavení (Open MPI) wrapperu: proměnná **OMPI\_FC** pro překladač, **OMPI\_FCFLAGS** pro volby překladače

### Loader mpirun/mpiexec (v Open MPI totéž co orterun)

Spouští požadovaný počet procesů na specifikovaných strojích.

- spuštění 4 procesů na jednom stroji  
**mpirun -np 4 a.out**
- spuštění procesů na více strojích  
rozpis slotů na strojích/uzlech v souboru **hostfile**:  
**slots** pro doporučený počet procesů, **max\_slots** pro maximální počet procesů  
počet procesorů počítače: př. **cat /proc/cpuinfo**  
při obsahu souboru **hosts**:  
node0 slots=2 max\_slots=8  
node1 slots=2 max\_slots=8  
po příkazu **mpirun --hostfile hosts -np 8 a.out** (volba **--byslot** je default)  
poběží na stroji node0 procesy 0,1,4,5 a na stroji node1 procesy 2,3,6,7 (cyklicky po slotech)  
po příkazu **mpirun --hostfile hosts -np 8 --bynode a.out**  
poběží na stroji node0 procesy 0,2,4,6 a na stroji node1 procesy 1,3,5,7 (cyklicky po uzlech)
- běh procesů v **agresivním** nebo **degradovaném** režimu  
agresivní režim (menší ochota procesu vzdát se procesoru) při počtu procesů ≤ počtu slots  
degradovaný režim po překročení počtu slots, tj. při nadobjednávce (oversubscribing) procesů  
vynucování agresivního režimu (př. je-li více slots než procesorů) je nevhodné (klesne výkon)
- procesová **afinita** (**mpi\_affinity**) pro přesné přiřazení procesů slotům pomocí souboru **rankfile** (volba **-rf rankfile**)  
př. **rankfile**: rank 0=node0 slot=0 // rank 1=node1 slot=0-3  
pozn. vláknová afinita není, ale vlákna procesu budou omezena na sloty předepsané procesu  
pozn. paměťová afinita smysluplná pouze na hardwaru s architekturou NUMA (non-uniform memory access)
- schéma **master-slave** (multiple programs multiple data, **MPMD**)  
**mpif90 -a.out master.f90 ; mpif90 -b.out slave.f90 ; mpirun -np 1 a.out : -np 3 b.out**

Ukázky: vzájemná výměna zpráv mezi dvěma procesy (ping-pong)

1. Blokující **MPI\_RECV**  
if (rank==0) then ; call MPI\_RECV ; call MPI\_SEND ; else ; call MPI\_RECV ; call MPI\_SEND ; endif  
oba procesy čekají v MPI\_RECV na příchod zprávy, která nepříjde: uváznutí (**deadlock**)
2. Synchronní **MPI\_SSEND**  
if (rank==0) then ; call MPI\_SSEND ; call MPI\_RECV ; else ; call MPI\_SSEND ; call MPI\_RECV ; endif  
oba procesy čekají v MPI\_SSEND na odezvu příjemce, která nepříjde: **deadlock**
3. Buferovaný **MPI\_BSEND**  
if (rank==0) then ; call MPI\_BSEND ; call MPI\_RECV ; else ; call MPI\_BSEND ; call MPI\_RECV ; endif  
procesy projdou přes MPI\_BSEND, pokud mají dostatečně velký bufer, jinak přetečení (**buffer overflow**)
4. Standardní **MPI\_SEND**  
if (rank==0) then ; call MPI\_SEND ; call MPI\_RECV ; else ; call MPI\_SEND ; call MPI\_RECV ; endif  
procesy projdou přes MPI\_SEND, pokud mají k dispozici (implicitní!) bufer, jinak deadlock synchronního sendu
5. **MPI\_SEND** a **MPI\_RECV** křížem, kombinovaný **MPI\_SENDRECV**
  - a) if (rank==0) then ; call MPI\_SEND ; call MPI\_RECV ; else ; call MPI\_RECV ; call MPI\_SEND ; endif  
MPI\_RECV čeká a dočká se, MPI\_SEND buď buferuje nebo čeká a dočká se: **spolehlivé řešení**
  - b) if (rank==0) then ; call MPI\_SENDRECV ; else ; call MPI\_SENDRECV ; endif  
MPI\_RECV čeká a dočká se, MPI\_SEND buď buferuje nebo čeká a dočká se: **spolehlivé řešení**
6. Ready mód **MPI\_RSEND**  
if (rank==0) then ; call MPI\_RSEND ; call MPI\_RECV ; else ; call MPI\_RECV ; call MPI\_RSEND ; endif  
MPI\_RSEND předpokládá připravený MPI\_RECV, ovšem (v Open MPI) není-li, MPI\_RSEND stejně počká
7. Neblokující **MPI\_ISEND**
  - a) if (rank==0) then ; call MPI\_ISEND ; call MPI\_RECV ; call MPI\_WAIT ;  
else ; call MPI\_ISEND ; call MPI\_RECV ; call MPI\_WAIT ; endif  
MPI\_ISEND nečeká, MPI\_RECV se dočká, MPI\_WAIT se také dočká: **spolehlivé řešení**
  - b) if (rank==0) then ; call MPI\_ISEND ; call MPI\_Irecv ; call MPI\_WAIT ; call MPI\_WAIT ;  
else ; call MPI\_ISEND ; call MPI\_Irecv ; call MPI\_WAIT ; call MPI\_WAIT ; endif  
MPI\_ISEND nečeká, MPI\_Irecv také nečeká, MPI\_WAIT se dočkají: **spolehlivé řešení**

Řešení č. 5 (send a recv křížem nebo kombinovaný sendrecv) je spolehlivé, může však být pomalejší než „duplexní“ řešení č. 3 s dostatečně velkým buferem (MPI\_ATTACH\_BUFFER). Řešení č. 7 je spolehlivé v obou variantách a rychlé alespoň jako č. 3, neboť nemusí nutně buferovat.