### Message Passing Interface

#### Nicolas Maillard

nicolas@inf.ufrgs.br

Instituto de Informática
Universidade Federal do Rio Grande do Sul

Curso de extensão Programação Paralela para Arquiteturas Multicores I.I. / INTEL





### Bibliografia

- http://www.mpi-forum.org
- Gropp, William et al., Using MPI, MIT Press.
- Gropp, William et al., Using MPI-2, MIT Press.
- Snir, M. et al., Dongarra, J., MPI: The Complete Reference.





## Interface de Passagem de Mensagens

#### Message Passing Interface

- Histórico:
  - PVM (Parallel Virtual Machine)
  - 1995 MPI 1.1
  - 1997 MPI 1.2
  - 1998 MPI-2
- 2 principais distribuições livres: LAM-MPI (www.lam-mpi.org) e MPI-CH (www-unix.mcs.anl.gov/mpi/mpich/).
- Distribuições de vendedores.
- "MPI is as simple as using 6 functions and as complicated as a user wishes to make it."





### Filosofia de MPI

- Single Program Multiple Data;
  - Vários processos executam todos o mesmo fluxo de instruções;
  - Cada um dos processos é identificado por seu rank.
- Troca de Mensagens:
  - existem primitivas especiais para comunicações "via rede";
  - o posto a um modelo de memória compartilhada.
- Acrescentar linguagens squenciais via biblioteca
  - Novos tipos de dados (MPI\_Status, MPI\_Communicator,...)
  - C, C++, Fortran.
  - Java (?)





# Seis instruções mágicas

- MPI\_Init Inicializa os processos.
- MPI\_Finalize Finaliza os processos.
- MPI\_Comm\_size determina o número de processos executando.
- MPI\_Comm\_rank determina o rank de um processo.
- MPI\_Send Manda uma mensagem.
- MPI\_Recv Recebe uma mensagem.





### MPI\_Init

- Declaração (em C): int MPI\_Init(int\*, char\*\*\*);
- Usado para inicializar os processos participando à execução, as estruturas de dados e para repassar os argumentos do main a todos os processos.
- Exemplo de chamada: MPI\_Init(&argc, &argv);
- Usa-se no início do programa!





### MPI\_Finalize

- Declaração (em C): int MPI\_Finalize();
- Libera os recursos, termina com o programa concorrente.
- Exemplo de chamada: MPI\_Finalize();
- Usa-se no fim do programa!





### MPI\_Comm\_size

- Declaração (em C): int MPI\_Comm\_size(MPI\_Communicator, int\*);
- Determina o número de processos em execução dentro do Communicator MPI.
- Exemplo de chamada:MPI\_Comm\_size(MPI\_COMM\_WORLD, &p).
- Usa-se em qualquer lugar do programa em geral no início.





### MPI\_Comm\_rank

- Declaração (em C): int MPI\_Comm\_rank(MPI\_Communicator, int\*);
- Determina o número (rank) do processo em execução dentro do Communicator MPI. O rank varia de 0 a p - 1.
- Exemplo de chamada:MPI\_Comm\_rank(MPI\_COMM\_WORLD, &r).
- Usa-se em qualquer lugar do programa em geral no início, para personalizar o comportamento do programa em função do rank.





## Meu primeiro programa MPI

- Veja /Ufrgs/Disciplinas/ProgParalelaPad/hello.c
- mpicc / mpirun no LAM.





### MPI\_Send

- int MPI\_Send(void\*, int, MPI\_Datatype, int, int, MPI\_Comm).
- Manda o conteúdo de um buffer do processo corrente para um processo destino.
- O buffer é determinado:
  - pelo tipo de dados no buffer (MPI\_Datatype);
  - por seu tamanho (número de itens no buffer);
- A mensagem é identificada por um tag.
- Exemplo de chamada: MPI\_Send(&work, 1, MPI\_INT, rank\_dest, WORKTAG, MPI\_COMM\_WORLD);





### MPI\_Recv

- int MPI\_Recv(void\*, int, MPI\_Datatype, int, int, MPI\_Comm, MPI\_Status\*).
- Recebe o conteúdo de uma mensagem em um buffer do processo corrente, vindo de um processo fonte.
- O buffer é determinado:
  - pelo tipo de dados no buffer (MPI\_Datatype);
  - por seu tamanho (número de itens no buffer);
- A mensagem é identificada por um tag.
- Existe um Status que possibilita a recuperação de informações sobre a mensagem após sua recepção.
- Exemplo de chamada: MPI\_Recv(&result, 1, MPI\_DOUBLE, 1, tag, MPI\_COMM\_WORLD, &status);





### Meu segundo programa MPI

- Veja /Ufrgs/Disciplinas/ProgParalelaPad/ping-pong.c
- Obs: só roda com  $p \ge 2!$  (PÉSSIMO EXEMPLO!)
- Obviamente, cada processo deve alocar/inicializar apenas as variáveis de que ele precisará, sem fazer o trabalho para as demais.





- MPI != Mestre/escravo.
- Esquema freqüente: testar a paridade do rank para diferenciar o comportamento.
- Os tags devem ser pre-definidos pelo programador e servem para evitar a colisão entre mensagens. Existe a constante MPI\_ANY\_TAG.
- O buffer é contíguo na memória. Se os dados a serem transmetidos estão numa estrutura espalhada, deve haver primeiro "serialização" dos mesmos. Cuidado com ponteiros...
- O buffer é de tamanho fixo. Caso se transmita valores em número desconhecido à compilação, deve-se usar 2 mensagens:
  - 1 primeira, de 1 int, para mandar o tamanho do buffer;
  - 1 segunda, contendo o buffer.





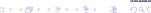
- MPI != Mestre/escravo.
- Esquema freqüente: testar a paridade do rank para diferenciar o comportamento.
- Os tags devem ser pre-definidos pelo programador e servem para evitar a colisão entre mensagens. Existe a constante MPI\_ANY\_TAG.
- O buffer é contíguo na memória. Se os dados a serem transmetidos estão numa estrutura espalhada, deve haver primeiro "serialização" dos mesmos. Cuidado com ponteiros...
- O buffer é de tamanho fixo. Caso se transmita valores em número desconhecido à compilação, deve-se usar 2 mensagens:
  - 1 primeira, de 1 int, para mandar o tamanho do buffer;
  - 1 segunda, contendo o buffer.





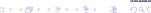
- MPI != Mestre/escravo.
- Esquema freqüente: testar a paridade do rank para diferenciar o comportamento.
- Os tags devem ser pre-definidos pelo programador e servem para evitar a colisão entre mensagens. Existe a constante MPI\_ANY\_TAG.
- O buffer é contíguo na memória. Se os dados a serem transmetidos estão numa estrutura espalhada, deve haver primeiro "serialização" dos mesmos. Cuidado com ponteiros...
- O buffer é de tamanho fixo. Caso se transmita valores em número desconhecido à compilação, deve-se usar 2 mensagens:
  - 1 primeira, de 1 int, para mandar o tamanho do buffer;
  - 1 segunda, contendo o buffer.





- MPI != Mestre/escravo.
- Esquema freqüente: testar a paridade do rank para diferenciar o comportamento.
- Os tags devem ser pre-definidos pelo programador e servem para evitar a colisão entre mensagens. Existe a constante MPI\_ANY\_TAG.
- O buffer é contíguo na memória. Se os dados a serem transmetidos estão numa estrutura espalhada, deve haver primeiro "serialização" dos mesmos. Cuidado com ponteiros...
- O buffer é de tamanho fixo. Caso se transmita valores em número desconhecido à compilação, deve-se usar 2 mensagens:
  - 1 primeira, de 1 int, para mandar o tamanho do buffer;
  - 1 segunda, contendo o buffer.





- MPI != Mestre/escravo.
- Esquema freqüente: testar a paridade do rank para diferenciar o comportamento.
- Os tags devem ser pre-definidos pelo programador e servem para evitar a colisão entre mensagens. Existe a constante MPI\_ANY\_TAG.
- O buffer é contíguo na memória. Se os dados a serem transmetidos estão numa estrutura espalhada, deve haver primeiro "serialização" dos mesmos. Cuidado com ponteiros...
- O buffer é de tamanho fixo. Caso se transmita valores em número desconhecido à compilação, deve-se usar 2 mensagens:
  - 1 primeira, de 1 int, para mandar o tamanho do buffer;
  - 1 segunda, contendo o buffer.





- No MPI\_Recv, pode-se usar o valor MPI\_ANY\_SOURCE como rank do processo de envio. Pode-se receber alguma-coisa de qualquer um, com qualquer tag! Neste caso, usa-se o MPI\_Status para recuperar as informações:
  - stat.MPI\_TAG
  - stat.MPI\_SOURCE
- O MPI\_Recv é bloqueante: o processo vai ficar esperando até ter recebido alguma coisa. Cuidado com Deadlocks.
- o MPI\_Send, logicamente, é bloqueante. Nas implementações, em geral ele não é (vide exemplo).





## Comunicações coletivas — definições

- São todas as comunicações que implicam mais de dois processos.
  - Oposto às comunicações "ponto-a-ponto".
- Sua prototipação no MPI permite a otimização, na biblioteca, da comunicação.
  - Nem sempre a implementação é tão eficiente assim. . .
  - Depende (a) da instalação; (b) da arquitetura.
- Exemplos: broadcast (one-to-all), all-to-all, all-to-one (reduce), split...
- Lembra muito do HPF!





### Sincronização global: barreira

- Pode ser preciso sincronizar todos (ou parte de) os processos.
  - Garantir que x processos passaram em um dado ponto do programa;
  - Medir um tempo.
- MPI\_Barrier(MPI\_Comm) MPI\_Barrier(MPI\_COMM\_WORLD);
- Obviamente, isso leva a perda de tempo!





### Broadcast (difusão)

- Um processo "mestre" difunde uma mensagem para x outros.
  - Vide aulas passadas para vários algoritmos (árvore binária, árvore de Fibonacci...);
  - obs: nada obriga a mandar para todos os p − 1 outros processos!
- MPI\_Bcast(void\*, int, MPI\_Datatype, int, MPI\_Comm)
- Manda os dados contidos no buffer, de tipo 'datatype', a partir do processo 'mestre', para todos os processos.

### Exemplo: Bcast de 1 int e de n doubles

```
MPI_Bcast(&size, 1, MPI_INT, 0, MPI_COMM_WORLD); #define RAIZ 0
```

MPI\_Bcast(dados, n, MPI\_DOUBLE, RAIZ, MPI\_COMM\_WORLD);

- Observações:
  - Todos os processos chamam MPI\_Bcast!
  - É uma sincronização global.
  - O buffer muda de semântica conforme o rank



#### Juntar dados

- Cada processo, inclusive o "mestre", manda dados para o mestre, que os armazena na ordem dos ranks.
- Exemplo típico de uma facilidade: pode-se fazer a mesma coisa com x MPI\_Send/MPI\_Recv!
- MPI\_Gather(void\*, int, MPI\_Datatype, void\*, int, MPI\_Datatype, int, MPI\_Comm)
- Argumentos: buffer de emissão, buffer de recepção, rank do mestre, comunicador.
- O argumento recvcount é o número de elementos recebidos / mandados por processo.

#### Exemplo: Gather de 10 doubles

```
double* sbuf, rbuf;
int mestre = 2;
if (rank == mestre) rbuf = malloc(p*10*sizeof(double));
else sbuf = malloc(10*sizeof(double));
MPI_Gather(sbuf, 10, MPI_DOUBLE, rbuf, 10, MPI_DOUBLE,
mestre, MPI_COMM_WORLD);
```



## Variações sobre as comunicações globais

- MPI\_Gatherv(...);
- MPI\_Scatter(...);
  - Faz o contrário do Gather!
  - Argumentos: buffer de emissão, buffer de recepção, rank do mestre, comunicador.
  - Existe também o Scatterv.
- MPI\_Allgather
  - Gather + cópia em todos os processos (Gather-to-all).
  - Argumentos: buffer de emissão, buffer de recepção, SEM mestre, comunicador.





#### All-to-All

- Versão mais "avançada" do Allgather, com dados distintos.
- O bloco j mandado pelo processo i vai ser recebido pelo processo j e armazenado no bloco i de seu buffer.
- MPI\_Alltoall(void\*, int, MPI\_Datatype, void\*, int, MPI\_Datatype, MPI\_Comm)
- Argumentos: buffer de emissão, de recepção, comunicador.

#### Exemplo: MPI\_Alltoall de 10 doubles

```
int i,j;
double* sbuf, *rbuf;
rbuf = malloc(p*10*sizeof(double));
sbuf = malloc(p*10*sizeof(double));
MPI_Alltoall(sbuf, 10, MPI_DOUBLE, rbuf, 10, MPI_DOUBLE,
MPI_COMM_WORLD); for (i=0; i<p; i++)
printf("Recebi de %d:", i);
for (j=0; j<10; j++) printf("%lf", rbuf[i*10+j]);
```



### Reduções

- Para efetuar uma operação comutativa e associativa em paralelo.
- MPI\_Reduce(void\*, void\*, int, MPI\_Datatype, MPI\_Op, int, MPI\_Comm)
- Argumentos: buffer de emissão, buffer de recepção, comunicador, tipo de operador.
- O operador pode ser MPI\_MAX, MPI\_MIN, MPI\_SUM, MPI\_PROD, MPI\_LAND, MPI\_LOR, etc...

#### Exemplo

```
double* sum_buf, resultado;
rbuf = malloc(p*10*sizeof(double));
sum_buf = malloc(10*sizeof(double));
MPI_Reduce(sum_buf, resultado, 10, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
```





## Comunicações coletivas: conclusões

- Procedimentos extremamente poderosos;
- Observa-se que s\u00e3o parecidos aos mecanimos de paralelismo de dados (HPF) ou de la\u00e7os (OpenMP): cf. o Reduce, por exemplo;
- Grande progresso rumo à especificação portável e eficiente do paralelismo;
- Obs: é um dos exemplos onde o MPI pode ser visto como modelo de programação mais avançado (!= soquetes, RMI, threads, ...).

### Comunicações ponto-a-ponto avançadas

- As operações básicas MPI\_Send e MPI\_Recv são bloqueantes e não locais:
  - não retornam até o usuário poder re-usar o buffer,
  - pode-se usar buffers intermediários (ou não);
  - o resultado do send depende do comportamento do processo que efetua o receive.
- Para deixar o usuário controlar exatamente o comportamento do programa, MPI oferece 3 modos de comunicação:
  - buffered: obriga as comunicações a usarem buffers, com controle de seu tamanho e uso, bem como do fim da operação;
  - síncrono: usa um Rendez-vous e retorna somente quando o Receive tem sido efetuado;
  - ready: o Send começa somente se o Recv tem começado; se não, retorna com erro.
- Existe só um MPI\_Recv, além de MPI\_Bsend, MPI\_Ssend e MPI\_Rsend.



### Vários MPI\_Send

- MPI\_Bsend, MPI\_Ssend e MPI\_Rsend têm o mesmo perfil como o MPI\_Send clássico;
- No caso do Bsend, MPI provê mais primitivas para controlar o uso dos buffers intermediários:
  - MPI\_Buffer\_attach(void\* buff, int size): define o espaço apontado por buff como sendo um buffer intermediário;
  - MPI\_Buffer\_detach(void\* buff, int size): libera o espaço apontado por buff de ser um buffer intermediário; essa chamada e bloqueante.





## Comunicações não-bloqueantes

- São mecanismos cruciais para obter alto-desempenho, a fim de poder mascarar as comunicações por cálculos;
- MPI usa uma estrutura de dados chamada MPI\_Request para identificar as comunicações não bloqueantes e poder verificar seu andamento;
- O perfil das chamadas é igual ao do MPI\_Send (ou MPI\_Recv), mas com um parâmetro final extra de tipo MPI\_Request.
- Pode-se mesclar o caráter não-bloqueante com os modos buffered, síncronos ou prontos.
- Exemplo:
   MPI\_request\* req; MPI\_Isend(&work, 1, MPI\_INT,
   rank\_dest, WORKTAG, MPI\_COMM\_WORLD, req);
   MPI\_Irecv(&result, 1, MPI\_DOUBLE, 1, tag,
   MPI\_COMM\_WORLD, req);
- Obs: as chamadas Irecv e ISend allocam e inicializam a Request.



## Controle do andamento das comunicações

- É preciso poder testar a conclusão das comunicações não bloqueantes para poder re-aproveitar seus buffers.
- Duas possibilidades:
  - MPI\_Wait(MPI\_Request\* req, MPI\_Status\* stat): retorna apenas quando a primitiva n\u00e3o bloqueante identificada por 'req' terminou. Ao retornar, 'stat' cont\u00e9m as informa\u00f3\u00f3es relevantes. 'req' \u00e9 liberada pelo Wait.
  - MPI\_Test(MPI\_Request\* req, int\* flag, MPI\_Status\* stat): testa se a comunicação terminou, sem bloquear. Ao retornar, 'flag' é setado e pode ser testado depois. Se o resultado é positivo, a 'req' é liberada e o 'stat' setado.
- Outras primitivas: MPI\_Request\_free; MPI\_Wait\_any; MPI\_Test\_any (com tabelas de requests).





## Controle do andamento das comunicações

- É preciso poder testar a conclusão das comunicações não bloqueantes para poder re-aproveitar seus buffers.
- Duas possibilidades:
  - MPI\_Wait(MPI\_Request\* req, MPI\_Status\* stat): retorna apenas quando a primitiva n\u00e3o bloqueante identificada por 'req' terminou. Ao retornar, 'stat' cont\u00e9m as informa\u00f3\u00f3es relevantes. 'req' \u00e9 liberada pelo Wait.
  - MPI\_Test(MPI\_Request\* req, int\* flag, MPI\_Status\* stat): testa se a comunicação terminou, sem bloquear. Ao retornar, 'flag' é setado e pode ser testado depois. Se o resultado é positivo, a 'req' é liberada e o 'stat' setado.
- Outras primitivas: MPI\_Request\_free; MPI\_Wait\_any; MPI\_Test\_any (com tabelas de requests).





## Controle do andamento das comunicações

- É preciso poder testar a conclusão das comunicações não bloqueantes para poder re-aproveitar seus buffers.
- Duas possibilidades:
  - MPI\_Wait(MPI\_Request\* req, MPI\_Status\* stat): retorna apenas quando a primitiva n\u00e3o bloqueante identificada por 'req' terminou. Ao retornar, 'stat' cont\u00e9m as informa\u00f3\u00f3es relevantes. 'req' \u00e9 liberada pelo Wait.
  - MPI\_Test(MPI\_Request\* req, int\* flag, MPI\_Status\* stat): testa se a comunicação terminou, sem bloquear. Ao retornar, 'flag' é setado e pode ser testado depois. Se o resultado é positivo, a 'req' é liberada e o 'stat' setado.
- Outras primitivas: MPI\_Request\_free; MPI\_Wait\_any; MPI\_Test\_any (com tabelas de requests).





### Estruturação dos processos

- MPI provê uma definição abstrata de conjuntos de processos:
  - para encapsular uma biblioteca MPI (espaço de nomes);
  - para possibilitar sincronizações de parte dos processos.
- MPI define grupos, contextos, topologia e comunicadores:
  - Grupo: coleção ordenada de processos, cada um tendo seu rank local ao grupo.
    - Define as comunicações ponto-a-ponto!
    - Define os particpantes em comunicações coletivas.
    - Default: MPI\_GROUP\_EMPTY.
  - Contextos: espaço de nomes para mensagens;
  - Topologia: organização virtual dos processos.
  - Comunicador: estrutura de dados que encapsula todo o resto.
    - Intra-comunicador: comunicação interna a um grupo;
    - Inter-comunicador: comunicação entre grupos;
    - Defaults: MPI\_COMM\_WORLD, MPI\_COMM\_SELF.
  - Obs: MPI não define o relacionamento rank/hostname.



### Estruturação dos processos

- MPI provê uma definição abstrata de conjuntos de processos:
  - para encapsular uma biblioteca MPI (espaço de nomes);
  - para possibilitar sincronizações de parte dos processos.
- MPI define grupos, contextos, topologia e comunicadores:
  - Grupo: coleção ordenada de processos, cada um tendo seu rank local ao grupo.
    - Define as comunicações ponto-a-ponto!
    - Define os particpantes em comunicações coletivas.
    - Default: MPI\_GROUP\_EMPTY.
  - Contextos: espaço de nomes para mensagens;
  - Topologia: organização virtual dos processos.
  - Comunicador: estrutura de dados que encapsula todo o resto.
    - Intra-comunicador: comunicação interna a um grupo;
    - Inter-comunicador: comunicação entre grupos;
    - Defaults: MPI\_COMM\_WORLD, MPI\_COMM\_SELF.
  - Obs: MPI n\u00e3o define o relacionamento rank/hostname.



### Operações com grupos

- MPI\_Group\_size(MPI\_Group g, int\* size) retorna o tamanho de um grupo.
- MPI\_Group\_rank(MPI\_Group g, int\* rank) retorna o rank dentro de um grupo.
- MPI\_Group\_translate\_ranks(MPI\_Group g1, int n, int\* ranks1, MPI\_Group g2, int\* ranks2) traduz os ranks de um grupo para o outro.
- MPI\_Comm\_group(MPI\_Comm comm, MPI\_Group\* group) retorna o grupo associado ao Comunicador 'comm'. Usado com o MPI\_COMM\_World!
- União, interseção, inclusão, ... de grupos são possíveis.
- MPI\_Group\_free(MPI\_Group g) libera o recurso.





### Topologia virtual com MPI

- simplifica a programação em alguns casos;
- exemplos: cartesiana; grafo;
- a topologia vem embutida no comunicador!
- MPI\_Cart\_create(MPI\_Comm old\_comm, int nbdims, int\* dims, int\* periodico, int reorder, MPI\_Comm\* comm\_cart);
- reorder: para re-ordenar os processos dentro da nova topologia.
- periodico[i] indica se a dimensão i é periodicamente mapeada.





### Topologia virtual com MPI

- simplifica a programação em alguns casos;
- exemplos: cartesiana; grafo;
- a topologia vem embutida no comunicador!
- MPI\_Cart\_create(MPI\_Comm old\_comm, int nbdims, int\* dims, int\* periodico, int reorder, MPI\_Comm\* comm\_cart);
- reorder: para re-ordenar os processos dentro da nova topologia.
- periodico[i] indica se a dimensão i é periodicamente mapeada.





### Operações com comunicadores

- MPI\_Comm\_size(MPI\_Comm c, int\* size) retorna o tamanho de um comunicador.
- MPI Comm rank(MPI Comm c, int\* rank) retorna o rank dentro de um comunicador.
- MPI\_Comm\_dup(MPI\_Comm c, MPI\_Comm new\_comm) efetua uma copia profunda da estrutura c.
- MPI\_Comm\_create(MPI\_Comm c, MPI\_Group g, MPI\_Comm\* new\_comm) cria um novo comunicador a partir do grupo g.
- MPI Comm\_split(MPI\_Comm c, int cor, int chave, MPI\_Comm\* new\_comm) cria um novo comunicador, através da partição de c, conforme for o valor de 'cor'. 'chave' permite definir o rank do processo em seu novo comunicador.
  - Útil para computação "Divisão & Conquista".
- MPI\_Comm\_free(MPI\_Comm c) libera o comunicador.



## MPI\_Datatypes

- MPI define tipos básicos (MPI\_INT, MPI\_DOUBLE, ...)
- MPI\_Byte (8 bits)
- MPI\_PACKED. Usado junto com MPI\_Pack(...) e MPI\_Unpack(...):
  - MPI\_Pack(void\* inbuf, int count, MPI\_Datatype dtype, void\* outbuf, int outcount, int\* position, MPI\_Comm comm)
    - Empacota em 'inbuf' os dados, atualiza "position" e retorna "outbuf".
  - MPI\_Send( outbuf, MPI\_Pack\_size( outbuf ), MPI\_PACK, ...)
  - MPI\_Unpack(...).





### Datatypes complexos

- O usuário pode definir um tipo complexo para MPI:
  - Define um conjunto de tipos básicos;
  - Define o deslocamento de cada tipo básico.
  - $\{(T_0, D_0), (T_1, D_1), \ldots, (T_n, D_n)\}.$
- A assinatura define a sequência de tipos usados;
- MPI\_Type\_struct( count, array\_of\_block\_length, array\_of\_deslocamento, array\_of\_types, new\_type);





### MPI-2

- "Nova" norma de MPI que inclui dinamicidade.
  - data de 1998!
  - (muito) recentemente implementada.
- Inclui:
  - criação dinâmica de processos;
  - comunicação entre processos dinamicamente criados;
  - RMA;
  - E/S paralelas.





### MPI\_Comm\_spawn

- Possibilidade de criar processos durante a execução do programa;
- Possibilidade de estabalecer comunicações entre 2 programas MPI rodando em paralelo.
- Usa extensivamente os comunicadores e cria um inter-comunicador para comunicação entre o grupo criador e o grupo criado.
- MPI\_Comm\_spawn( char\* comando, char\* argumentos[], int maxprocs, MPI\_Info info, int root, MPI\_Comm com, MPI\_Comm\* inter\_comm, int erros[]);
- dispara 'maxprocs' processos para executar o 'comando' (= outro programa MPI). Os argumentos devem ter sentido no processo 'root'. Na saída, 'inter\_comm' contém um comunicator para possibilitar a comunicação entre os novos processos e os antigos.

### Comunicação entre processos em MPI-2

- O mecanismo básico fica igual: troca de mensagens.
- Mensagens ponto-a-ponto: usa-se o inter-comunicator retornado pelo MPI\_Comm\_Spawn.
  - A maior diferença é que o processo fonte pode acessar ranks de processos dependendo do número no segundo comunicador!
- Mensagens coletivas: é uma novidade do MPI-2. Funciona da mesma forma.
- Pode-se também usar MPI\_Intercomm\_merge para fusionar os comunicadores em um só, e se comunicar normalmente nele depois.





### Conectar aplicações MPI-2

- Pode-se usar mecanismos tipo cliente/servidor entre (grupos de) processos MPI-2.
- Lado servidor: Criação de porta e aceitação de conexão a ela;
  - MPI\_Open\_port(MPI\_Info, char\* nome-porta);
  - MPI\_Comm\_accept(char\* nome-porta, MPI\_Info, int root, MPI\_Commm, MPI\_Comm\* inter\_com) — chamada coletiva e bloqueante que retorna um inter-comunicador.
- Lado cliente: Conexão a uma porta.
  - MPI\_comm\_connect(char\* nome-porta, MPI\_Info info, int root, MPI\_Comm comm, MPI\_Comm\* inter\_comm);
- uma vez que foi estabelecida a conexão, pode-se usar os Send/Recv clássicos através do inter-comunicator.
- Obs: existe um MPI\_Publish\_name e um MPI\_Lookup\_name.





### Acesso distante a memória (RMA)

- Alternativa à troca de mensagens!!!
- Outro modelo de programação paralela
  - Escritas/leituras diretas, remotas, na memória de um outro processo (sem sua cooperação).
- MPI\_Win\_create para definir uma janela de memória
- MPI\_Put, MPI\_Get para escrever e ler. accessível.





### E/S paralelas

- Solução clássica com MPI-1.2: um processo faz as E/S para/de um arquivo...
- Melhora: cada processo escreve para um arquivo diferente.
- Solução com MPI-2:
  - MPI\_File em lugar de FILE;
  - MPI\_File\_open MPI\_File\_write MPI\_File\_read MPI\_File\_close
  - usa um comunicator para definir quais processos acessam ao arquivo;
  - escreve o conteúdo de um buffer no arquivo;
  - Usa MPI\_File\_view para especificar o deslocamento no arquivo onde cada processo vai escrever.
    - MPI\_File\_set\_view(arquivo, rank\*BUFSIZE\*sizeof(char), MPI\_CHAR,...)



