

Progra° en mémoire distribuée I

MPI: message passing interface

MPI: librairie de f° de comm° pr l'envoi & récep° de msg entre processus. SPMD. regroupe ts les processus  
 Communicateur: regroupe les processus qui peuvent échanger des msg. MPI\_COMM\_WORLD: 3 par défaut  
 Types de communi°: pt par pt ou collectives, doivent être au m° comm

Comm° synchrone: dmde de transmission, accepta°, envoi de données VS asynchrone: émetteur sait seul que msg envoyé  
 Comm° non bloquante: on vérif + tard si comm° réussie. Bloquante: on attend. Synchrone: f° bloquante.  
 Com° avec buffer: on fait une copie du msg, émetteur non bloqué si récepteur non disp°.  
 Com° collectives & Barrières: pgr arrêt until ts proc arrivés à la barrière  
 ↳ Réduction: com° 1 à tous ↳ Réduction: donnée unit° résultat de ≠ proc.

#include <mpi.h> Processus liés et n°és from 0.  
 MPI: Init (int \*argc, char \*\*argv[]); 1° f° à exécuté, crée MPI\_COMM\_WORLD. Fin: MPI: Finalize();  
 int MPI: Comm: Get\_size(): retourne le nb de processus du communicateur Comm  
 int MPI: Comm: Get\_rank(): retourne identificateur proc ds communicateur Comm

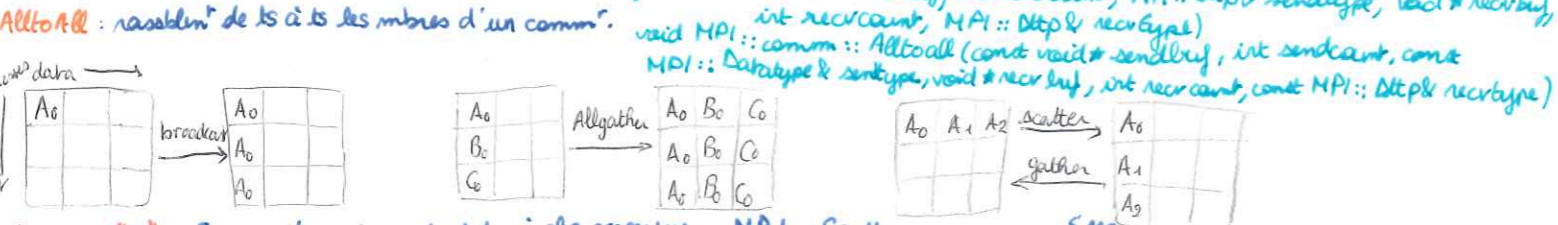
F° de comm°: send et receive. Possibilité de deadlock. Spécifier récepteur/émetteur, données (début, fin), type, identificateur de msg  
 MPI: Comm: Send (const void \*buf, int count, MPI: Datatype & datatype, int dest, int tag) Datatypes: MPI\_CHAR, (tag) MPI\_FLOAT...  
 MPI: Comm: Recv (void \*buf, int count, MPI: Datatype & datatype, int source, int tag, MPI: Status & status)

Obt taille msg: Status: Get\_count (const Datatype & datatype)  
 Pr Recv: émetteur peut être MPI\_ANY\_SOURCE, tag peut être MPI\_ANY\_TAG  
 MPI: Status \* status a comme info: MPI: Status: Get\_source(): Émetteur msg  
 int MPI: Status: Get\_tag(): tag du msg reçu  
 int MPI: Status: Get\_error(): code d'erreur

Mesurer temps: double MPI: Wtime(). T1 = MPI: Wtime(); Tmp = T2 - T1;  
 PMD II Comm° collective: bloquante, ts les proc du comm° doivent exécuter la f° collective, buf récep° doit avoir l'exacte msg

Types de mvmt de données

Radia°: Un proc racine send data au reste: void MPI: Comm: Bcast (void \*buff, int count, const MPI: Datatype & dt, int root)  
 Scatter: Distrib° data d'1 proc aux autres: void MPI: Comm: Scatter (void \*buff, sendbuf, int sendcount, MPI: Dtp & sendtype, void \*recvbuf, int recvcount, MPI: Dtp & recvtype, int root)  
 Gather: Rassembl° data de ts proc ds 1  
 void MPI: Comm: Gather (const void \*sendbuf, int sendcount, const MPI: Dtp & sendtype, void \*recvbuf, int recvcount, const MPI: Dtp & recvtype, int root)  
 AllGather: gather de ts à ts. void MPI: Comm: Allgather (const void \*sendbuf, int sendcount, MPI: Dtp & sendtype, void \*recvbuf, int recvcount, MPI: Dtp & recvtype, int root)  
 AlltoAll: rassembl° de ts à ts les mbres d'un comm°.  
 void MPI: Comm: Alltoall (const void \*sendbuf, int sendcount, const MPI: Datatype & sendtype, void \*recvbuf, int recvcount, const MPI: Dtp & recvtype)



Variants "v": Envoi d'un nb de data à chq processus. MPI: Scatter, MPI: Gather, MPI: Allgather, MPI: Alltoall  
 Calcul en groupe: # d'elts/entrées à chq proc déf par des vecteurs de P elts, P: # de proc du communicateur.

Réduction: op° avec data de chq proc, rés dans root. void MPI: Comm: Reduce (const void \*sendbuf, void \*recvbuf, int count, const MPI: Dtp & dt, const MPI: Op & op, int root).  
 Op° de réduc° prédéfinies: MPI: [MAX, MIN, SUM, PROD, LAND, LOR, BAND, BOR, LXOR, BXOR, MAXLOC, MINLOC].  
 Synchr°: Barrière: arrêt until ts arrivés: int MPI: Barrier (MPI: Comm comm)

PMD III Comm° en arrière: int MPI: sendrecv(...) contient params de MPI: send et recv ensemble, énte deadlock  
 Send/Rcv non bloquant: int MPI: Irecv (... , MPI: Request \* request) copie msg ds buffer et retourne tds  
 Attend finalisa° comm°: int MPI: Wait (MPI: Request \* request, MPI: Status \* status)  
 int MPI: Test ("request", int \*flag, "status"): flag = 0 si comm non complète

Tester si msg arrivent sans attendre, connaître émetteur, libelle, taille: MPI: Probe & MPI: Iprobe (avt Recv ou Irecv)  
 Avec buffer: MPI: Bsend (m° args q Send). À faire: MPI: Buffer - [detach] (void \*buffer, int size)  
 Récepteur prêt immédiatement: MPI: Rsend (ready)  
 Envoi synchrone: MPI: Ssend (retourne seul qd msg reçu et commence à lire)

Versions non bloquantes: MPI: [Irecv, Irecv, Bsend, Ssend, Ssend]  
 Comm° persistante: m° comm° gd nb de x, m° args (vals <>)  
 ↳ Création du contexte: int MPI: [Send, Recv] - [attach] (... , MPI: Request \* request)  
 ↳ Envoi du message: int MPI: Start (MPI: Request \* request)  
 ↳ Astuces: Irecv: besoin de comm° non bloquantes, Bsend et Rsend: cas très particuliers

Irecv: la + habituelle, Ssend: best result car pas buffers

PROGRA //

#include <omp.h>

HPC: High performance computing (algos, archi, progra //)  
 SPMD: Single Program Multiple Data → mémoire partagée: OpenMP (MultiProcessing) pr C, C++, Fortran  
 → mémoire distribuée: MPI

Progra° en mémoire partagée I: La coordina° et coop° st faites en usant R/W de vars partagées + vars de synchr°.  
 → Bas niveau: barrières, locks, sec° critiq, sémaphores...  
 → Haut niveau: directives de compila° (OpenMP).

Construct° de régions // parallel  
 #pragma omp parallel {...} int id = omp\_get\_thread\_num();  
 Def #threads: (void) omp\_set\_num\_threads(8); Nb threads en exec°: int omp\_get\_num\_threads();  
 #pragma omp parallel num\_threads(8) {...}

int main (int argc, char \*argv[]) { ... int a = atoi(argv[1]);  
 clauses dans une région //:  
 #pragma omp parallel if (a > 0) {...}  
 #pragma omp parallel private (a) shared (b, c): vars non init  
 #pragma omp parallel firstprivate (a): var init mais sont pas  
 #pragma omp parallel reduction (+:b): res final de + \* max min and or

PMPI II Partage du travail: for, sections, single  
 #pragma omp (//) for: barrière implicite à la fin de la bcle à - que newait  
 ↳ private ↳ firstprivate ↳ reduction ↳ lastprivate: val de la dernière it°/thread  
 #pragma omp for schedule (type [, taille bloc]): poliq de partage des it° d'1 bcle parmi threads  
 ↳ static: it° % en bloc de taille chunk assignés en round-robin, default = num\_it / num\_threads  
 ↳ dynamic: 1° th fini, 1° serv. default = 1. ↳ runtime: var d'env° OMP\_SCHEDULE  
 ↳ guided: taille chunk ≈ n it / n th et ds exp° ↳ auto: poliq laissée au compilateur ou au runtime

Astuce: static si it° équilibrées, static, n pr déséquilibre léger, dynamic: @ si charges <>, mais - localité data  
 ↳ guided: - coût q dynq mais (!) qd 1° it° les + chères  
 #pragma omp for collapse (val - pos): it° groupées pr former 1! espace d'it° réparti entre threads

#pragma omp [//] sections [clauses] { #pragma omp sec° [ bloc code ] #p. omp sec° [ bloc ... ]  
 ↳ private ↳ firstprivate ↳ lastprivate ↳ red° ↳ newait  
 #pragma omp single [clauses] { ... } : exec par 1! thread  
 ↳ private ↳ firstprivate ↳ newait ↳ copyprivate: pas avec newait. val (var privée) copiée aux autres threads à la fin.

PMD III Synchronisation: master, critical, atomic, barrier, ordered

#pragma omp master: Seul le th master exec bloc les autres sautent Plus de barrières implicites  
 #pragma omp critical [nom]: Ts threads exec bloc, mais 1! th a accès au code en m° temps  
 #pragma omp atomic: x++, ++x, x--, --x, x binop = exp, binop: + \* - / & ! ^ << >>  
 - coûteux que critical, ne s'applique qu'à la sentance code immédiatement sur Protège maj de var partagée

#pragma omp barrier: on attend que ts les threads soient arrivés à cette ligne  
 #pragma omp ordered: bloc exec ds le m° ordre qu' on séquentiel  
 Gestion de tâches: task, taskwait

#pragma omp task [clauses]: crée une new tâche (unité de travail) qui sera exec par le th l'ayant créée  
 ↳ firstprivate ↳ private ↳ shared ↳ default (private | shared | none) ↳ unbind: tâche plus liée à th donnée  
 ↳ if (expr): si True, th exec tâche, sinon différé

#pragma omp taskwait: on attend la terminaison de toutes les tâches