

Programmation parallèle hybride

Marc Tajchman

CEA - DEN/DM2S/STMF/LMES

10/12/2020

Cas envisagés

Du fait qu'il y a souvent plusieurs dispositifs matériels (clusters, machines multi-cœurs, GPU, etc.) et plusieurs outils logiciels associés, il est intéressant d'essayer de combiner leur utilisation.

En fonction de ce qui est disponible sur les machines

1. clusters (plus généralement, machines parallèles multi-nœuds) : MPI,
2. processeurs multi-cœurs : OpenMP (et autres outils : TBB, std::threads, ...),
3. accélérateurs de calcul (GPU ou autres) : Cuda, OpenCL,
4. outils PGAS ("partitionned global addressing system") :
pour information (faibles performances),

on testera plusieurs combinaisons.

Machines multi-nœuds et multi-cœurs

Pour utiliser la puissance de calcul de ce type de machine, on a souvent le choix entre 2 possibilités :

- ▶ **MPI** pour le parallélisme multi-nœuds et **MPI** pour le parallélisme multi-cœurs (plusieurs processus MPI sur chaque nœud).
- ▶ **MPI** pour le parallélisme multi-nœuds et **OpenMP** pour le parallélisme multi-cœurs (un processus MPI sur chaque nœud).

Sur des simulations de taille petite ou moyenne, on préfère souvent le premier choix pour des raisons de simplicité.

Par contre, sur des simulations de (très) grande taille, le “tout MPI” atteint plus rapidement ses limites d'utilisation.

Avantages/inconvénients du tout MPI:

- ⊕ programmation MPI nœuds-cœurs identique à la programmation MPI entre les nœuds
- ⊕ les bibliothèques MPI sont de mieux en mieux optimisées (en particulier si plusieurs processus sont sur le même processeur)
- ⊖ le nombre de processus MPI augmente plus vite (n° cœurs \times n° nœuds), les structures internes de MPI, la mémoire utilisée pour les communications aussi (“cellules fantômes”, “halos”)

On atteint actuellement les limites des machines parallèles sur des simulations avec des nombres de processus $> 10^6$.

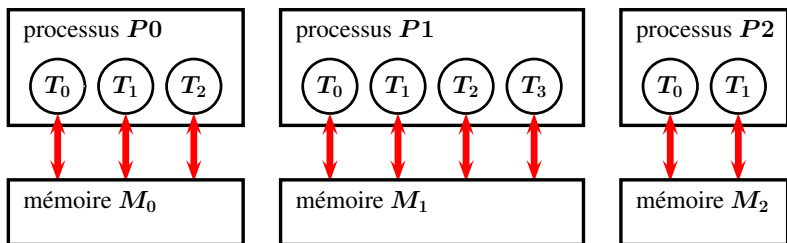
Avantages/inconvénients de la combinaison MPI-OpenMP:

- ⊕ les possibilités en nombre de nœuds-cœurs sont plus importantes
Le nombre total de processus MPI diminue (un seul processus MPI par nœud de calcul.)
- ⊖ la programmation MPI sur les nœuds et OpenMP sur les cœurs est plus complexe
- ⊖ l'amélioration des performances n'est pas toujours évidente surtout pour des nombres de processus petits ou moyens.

Exemples4/MemoireMPI : ressource mémoire utilisée par MPI (MPI_Init et MPI_Finalize, seulement, pas d'échange de messages).

Lire le fichier README.txt pour des instructions

Combinaison des modèles MPI et OpenMP en modèle hybride MPI-OpenMP:



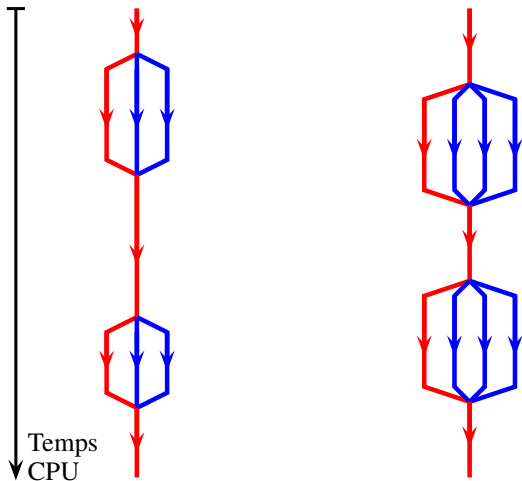
- ▶ Les différents processus P_0, P_1, \dots ont chacun leur mémoire locale M_0, M_1, \dots
- ▶ A l'intérieur de chaque processus P_i , un ou plusieurs threads travaillent avec la mémoire locale M_i .

- ▶ Pas de synchronisation a priori entre les sections multi-threads entre les différents processus, donc besoin éventuel de combiner des barrières OpenMP et MPI
- ▶ 2 types de réductions dans OpenMP et MPI, donc pour faire une réduction complète, il faut combiner ces 2 réductions
- ▶ On peut a priori échanger des messages entre chaque threads de 2 processus différents

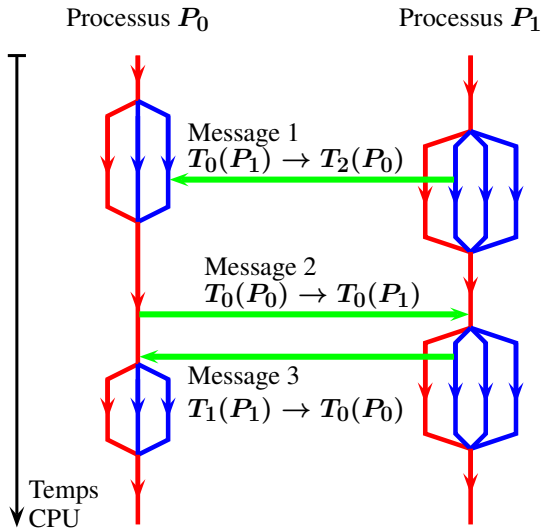
Exemple sur 2 processus, chacun avec plusieurs threads

Processus P_0

Processus P_1



Quelques types différents de messages



Le grand nombre de situations possibles est très difficile à gérer pour ceux qui développent les libraires OpenMP et MPI

Pour cette raison, 4 niveaux de compatibilité entre OpenMP et MPI, ont été définis:

- ▶ **MPI_THREAD_SINGLE:**
Un seul thread par processus (donc pas d'OpenMP)
- ▶ **MPI_THREAD_FUNNELED:**
Parmi les threads, seul celui qui a initialisé MPI peut utiliser des fonctions MPI
- ▶ **MPI_THREAD_SERIALIZED:**
A un instant donné, un seul thread peut utiliser des fonctions MPI (mais pas toujours le même thread)
- ▶ **MPI_THREAD_MULTIPLE:**
Pas de restrictions, tous les threads peuvent utiliser les fonctions MPI.

Plus le niveau de compatibilité permet de souplesse dans l'utilisation des fonctions MPI et des pragmas OpenMP, plus la gestion interne des threads dans les processus est coûteuse. **Donc, il faut choisir le niveau le plus bas suffisant pour l'algorithme du code**

Pour choisir le niveau de compatibilité, il faut remplacer

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

par

```
MPI_Init_thread(int *argc , char ***argv ,  
                int required , int *provided)
```

où `required` est un entier qui peut prendre une des 4 valeurs

- ▶ `MPI_THREAD_SINGLE`,
- ▶ `MPI_THREAD_FUNNELED`,
- ▶ `MPI_THREAD_SERIALIZED`,
- ▶ `MPI_THREAD_MULTIPLE`

et `provided` est un entier qui reçoit le niveau de compatibilité effectivement proposé par la version de MPI

Il est important de comparer les valeurs de `required` et `provided` pour vérifier que le niveau demandé est effectivement disponible.

Exemples4/Hybrides : différents niveaux de combinaison MPI-OpenMP