

Correction du Partiel de PP, 2019-2020

Polio, Promo 2021

23 novembre 2019

Durée : 1h30. Ce partiel est découpé en 5 exercices.

Correction par Polio (promo 2021), avec l'énoncé transcrit avant chaque réponse.

1 Exercice I : Petit programme MPI

Soit la section de code suivante exécutée par tous les processus d'un programme MPI (il y a au moins 2 processus) :

```
1 int rang, som;
2
3 MPI_Comm_rank(MPI_COMM_WORLD, &rang);
4
5 som = 0;
6 if (rang == 0)
7     som = 1;
8 som += 2;
9 if (rang == 1)
10    som += 3;
```

Suite à l'exécution de ces lignes, quelle est la valeur de la variable `som` en fonction du rang MPI (justifiez brièvement) ?

Tous les processus passent par `som = 0` et `som += 2`. Ainsi les processus dont le rang est supérieur ou égal à 2 auront pour valeur finale `som`: 2.

Rang 0 : Passe par `som = 1` après le `som = 0`, la valeur finale sera donc `som`: 3.

Rang 1 : Passe par `som += 3` après le `som += 2`, la valeur finale sera donc `som`: 5.

2 Exercice II : Cherchez l'erreur

Pour chacune des sections de code suivantes (numérotées **a** à **d**), dites :

- si elle fait l'objet d'un risque de blocage ("*deadlock*") ;
- ou bien si les contenus des messages risquent d'être erronés ;

Justifiez vos choix et apportez une correction pour que la section soit correcte.

Dans tout ce qui suit, on désigne par `rang` le rang du processus MPI et par `P` la taille du communicateur `MPI_COMM_WORLD`. On suppose que le programme a correctement initialisé MPI.

2.1 Section A

On veut calculer le nombre total d'éléments nuls dans un tableau distribué.

```
1 int i, Nnuls = 0, Nnuls_tot;
2 /* tabl est un tableau distribue sur tous les processus et prealablement initialise
3  * N = dimension locale du tableau
4  */
5 for (i = 0; i < N; i++)
6     if (tabl[i] == 0)
7         Nnuls++;
8
9 if (Nnuls > 0) {
10    MPI_Allreduce(&Nnuls, &Nnuls_tot, 1, MPI_INT, MPI_SUM MPI_COMM_WORLD);
11    printf("Nb d'elements nuls dans le tableau distribue : %d\n", Nnuls_tot);
12 }
```

`MPI_Allreduce` est une opération collective : tous les processus compris dans `MPI_COMM_WORLD` doivent l'appeler. Ici, seuls les processus ayant trouvé au moins un élément nul l'appellent, il y aura donc un *deadlock*.

Pour corriger, il faut supprimer le `if (Nnuls > 0)`.

2.2 Section B

Le processus 0 veut distribuer un entier différent à chacun des autres processus

```
1 MPI_Request req[P];
2 int entier;
3
4 if (rang == 0)
5 {
6     for (i = 1; i < P; i++)
7     {
8         entier = i * i;
9         MPI_Isend(&entier, 1, MPI_INT, i, 1000, MPI_COMM_WORLD, req + (i - 1));
10    }
11    MPI_Waitall(P - 1, req, MPI_STATUSES_IGNORE);
12 }
13 else
14 {
15     MPI_Recv(&entier, 1, MPI_INT, 0, 1000, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
16 }
```

La variable `entier` est utilisée comme stockage pour chaque `MPI_Isend(&entier, [...])`. Si un de ces envois n'est pas complété avant que `entier` ne soit réécrit par le tour de boucle suivant, la valeur envoyée sera alors erronée.

Une correction consiste à déclarer un tableau `int entiers[P-1]` et à remplacer le contenu de la boucle par :

```
1 entiers[i - 1] = i * i;
2 MPI_Isend(entiers + (i - 1), 1, MPI_INT, i, 1000, MPI_COMM_WORLD, req + (i - 1));
```

Il faut toutefois laisser la variable `entier` afin qu'elle puisse être utilisée pour la réception pour les autres processus.

2.3 Section C

But : faire circuler un message (le contenu de `m`) dans un anneau initialisé par le processus 0. Le message circule dans le sens croissants des rangs modulo `P`.

```
1 int droite, gauche, m;
2
3 if (rang == 0)
4 {
5     m = 0;
6 }
7
8 droite = (rang + 1) % P;
9 gauche = (rang - 1) % P;
10
11 MPI_Recv(&m, 1, MPI_INT, gauche, 1000, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
12 MPI_Send(&m, 1, MPI_INT, droite, 1000, MPI_COMM_WORLD);
```

`MPI_Recv` (et `MPI_Send`) est bloquante, hors tous les processus commencent par l'appeler : aucun `MPI_Send` ne sera fait. Il y a donc un *deadlock*.

Une correction possible consiste à spécialiser le code pour le processus 0 :

```
1 if (rang == 0)
2 {
3     MPI_Send(&m, 1, MPI_INT, droite, 1000, MPI_COMM_WORLD);
4     MPI_Recv(&m, 1, MPI_INT, gauche, 1000, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
5 }
6 else
7 {
8     MPI_Recv(&m, 1, MPI_INT, gauche, 1000, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
9     MPI_Send(&m, 1, MPI_INT, droite, 1000, MPI_COMM_WORLD);
10 }
```

2.4 Section D

Le processus 0 veut envoyer le message `m = 100` au processus 1.

```
1 if (rang == 0)
2 {
3     m = 100;
4     MPI_Request req;
5     MPI_Isend(&m, 1, MPI_INT, 1, 1000, MPI_COMM_WORLD, &req);
6 }
7 else if (rang == 1)
8 {
9     MPI_Recv(&m, 1, MPI_INT, 0, 1000, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
10 }
```

`MPI_Isend` est non bloquant : le message risque de ne jamais être envoyé et donc de bloquer le processus 1 indéfiniment.

Pour corriger il est possible d'ajouter la ligne `MPI_Wait(&req, MPI_STATUS_IGNORE);` en dessous de la ligne d'envoi.

3 Exercice III : Quelle est la bonne valeur ?

Soit la section de code suivante exécutée par tous les processus d'un programme MPI (il y a au moins 3 processus) :

```
1 int rang, val, valA, valB;
2 MPI_Request req[2];
3
4 MPI_Comm_rank(MPI_COMM_WORLD, &rang);
5
6 if (rang == 0)
7 {
8     MPI_Irecv(&valA, 1, MPI_INT, MPI_ANY_SOURCE, 1000, MPI_COMM_WORLD, req + 0);
9     MPI_Irecv(&valB, 1, MPI_INT, MPI_ANY_SOURCE, 1111, MPI_COMM_WORLD, req + 1);
10    MPI_Waitall(2, req, MPI_STATUSES_IGNORE);
11 }
12 else if (rang == 1)
13 {
14     val = rang;
15     MPI_Send(&val, 1, MPI_INT, 0, 1111, MPI_COMM_WORLD);
16 }
17 else if (rang == 2)
18 {
19     val = rang;
20     MPI_Send(&val, 1, MPI_INT, 0, 1000, MPI_COMM_WORLD);
21 }
```

Question : Suite à l'exécution de ces lignes, quelles sont les valeurs des variables `valA` et `valB` du processus de rang 0 (justifiez) ?

La différenciation sur les réceptions du processus 0 se fait sur les tags (1000 et 1111).

Le processus 1 envoie sur le tag 1111 donc `valB = 1`. Le processus 0 envoie sur le tag 1000 donc `valA = 2`.

4 Exercice IV : Implémentation d'un allgather

Voici le code de la fonction `allgather` et son programme principal : complétez-le.

```
1 void allgather(int in, int *out)
2 {
3     int rang, P;
4
5     /* 1. A COMPLETER: il faut affecter le nombre de processus dans 'P' */
6     /* 2. A COMPLETER: il faut affecter la variable 'rang' */
7
8     if (rang == 0)
9     {
10        MPI_Status sta;
11        out[0] = in;
12        for (int p = 1; p < P; p++)
13        {
14            int val;
15            /* 3. A COMPLETER: reception dans 'val' de la valeur envoyee par le processus de rang
16            'p' */
17            /* 4. A COMPLETER: mise a jour du tableau 'out' */
18        }
19    }
20    else
21    {
22        /* 5. A COMPLETER: envoi de 'in' au processus 0 */
23    }
24    /* 6. A COMPLETER: le processus 0 diffuse 'out' aux autres processus */
25 }
26
27 int main(int argc, char *argv[])
28 {
29     int rang, P, *out;
30
31     /* 7. A COMPLETER */
32     /* 8. A COMPLETER: il faut affecter le nombre de processus dans 'P' */
33     /* 9. A COMPLETER: il faut affecter la variable 'rang' */
34
35     out = (int *)malloc(P * sizeof(int));
36     allgather(rang, out);
37
38     for (int p = 0; p < P; p++)
39         printf("P%d, out[%d] = %d\n", rang, P, out[p]);
40     free(out);
41
42     /* 10. A COMPLETER */
43
44     return 0;
45 }
```

Soit la fonction `void allgather(int in, int *out)` appelée en même temps par tous les processus MPI. Elle collecte dans le tableau `out` le contenu de toutes les variables `in` de tous les processus MPI. Tous les processus doivent récupérer ensuite les résultats. Le tableau `out` est préalable alloué au nombre total de processus.

La fonction `allgather` est implémentée dans le bloc de code ci-dessus, ainsi qu'un programme principal qui utilise cette fonction. Cependant, ce code n'est pas complet.

Travail à faire : compléter les lignes `/* À COMPLÉTER */` pour que la fonction `allgather` et le programme principal soient corrects.

- 1 et 8 : `MPI_Comm_size(MPI_COMM_WORLD, &P);`
- 2 et 9 : `MPI_Comm_rank(MPI_COMM_WORLD, &rang);`
- 3 : `MPI_Recv(&val, 1, MPI_INT, p, 1000, MPI_COMM_WORLD, &sta);`
- 4 : `out[p] = val;`
- 5 : `MPI_Send(&in, 1, MPI_INT, 0, 1000, MPI_COMM_WORLD);`
- 6 : `MPI_Bcast(out, P, MPI_INT, 0, MPI_COMM_WORLD);`
- 7 : `MPI_Init(&argc, &argv);`
- 8 : `MPI_Finalize();`

5 Exercice V : Recouvrement

Définition de la fonction `collect_trait` :

```
1 void collect_trait(int *tab, int N) {
2     int rank, P, k;
3     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
4     MPI_Comm_size(MPI_COMM_WORLD, &P);
5
6     if (rank == 0) {
7         MPI_status sta;
8         int *buf2 = (int *)malloc(N * sizeof(int));
9         trait(tab, N);
10        for (k = 1; k < P; k++) {
11            MPI_Recv(buf2, N, MPI_INT, k, 1000, MPI_COMM_WORLD, &sta);
12            trait(buf2, N);
13        }
14        free(buf2);
15    }
16    else {
17        MPI_Send(tab, N, MPI_INT, 0, 1000, MPI_COMM_WORLD);
18    }
19 }
```

On considère une fonction `collect_trait` dont la définition se trouve ci-dessus. Elle est utilisée dans les conditions suivantes :

1. Cette fonction est appelée par tous les processus MPI (comme une opération collective);
2. Chaque processus a pré-alloué et rempli son propre tableau de `N` entiers;
3. La valeur `N` est identique pour tous les processus MPI;
4. La fonction `trait(int *tab, int N)` (appelée par le processus 0) effectue un traitement coûteux à partir des données du tableau `tab` passé en argument, sans modifier le contenu de ses données.

Question 1 : Quelle est la modification à apporter à la réception pour que le processus 0 reçoive un message à partir de n'importe quel processus ?

Il faut remplacer le `k` utilisé pour indiquer depuis quel processus réceptionner par `MPI_ANY_SOURCE`.

Question 2 : Expliquez brièvement pourquoi la réception (par le processus 0) n'est pas recouverte par le traitement `trait`.

`MPI_Recv` est bloquant : la réception est donc forcément effectuée avant que le traitement des données reçues ne s'effectue.

Question 3 : Modifiez la section de code exécutée par le processus 0 pour recouvrir la réception par le traitement `trait` tout en recevant les messages dans n'importe quel ordre.

Indication : utilisez deux buffers de `N` entiers, un pour recevoir les données, l'autre pour traiter les données déjà reçues.

```
1 MPI_status sta;
2 // Nouvelles variables utilisees pour le recouvrement
3 MPI_Request req;
4 int n;
5 int *buf2 = (int *)malloc(N * sizeof(int));
6 int *recvbuf = (int *)malloc(N * sizeof(int));
7 // Copie du tableau du rang 0 pour qu'il soit traite pendant la premiere reception
8 for (n = 0; n < P; n++) buf2[n] = tab[n];
9 // Recouvrement avec reception depuis toutes les sources
10 for (k = 1; k < P; k++) {
11     MPI_Irecv(recvbuf, N, MPI_INT, MPI_ANY_SOURCE, 1000, MPI_COMM_WORLD, &req);
12     // Traitement du tableau precedent pendant la reception du suivant
13     trait(buf2, N);
14     MPI_Wait(&req, &sta);
15     for (n = 0; n < P; n++) buf2[n] = recvbuf[n]; // Copie
16 }
17 // Traitement du dernier tableau qui n'a pas ete fait apres la derniere copie: il faut donc le
   faire apres la boucle
18 trait(buf2, N);
19 free(buf2);
20 free(recvbuf); // Ne pas oublier de liberer le nouveau buffer
```