



Informatique Parallèle et Distribuée

Notes de laboratoire

Pr. Manneback Pierre
Dr. Sidi Ahmed Mahmoudi
Ir. François Rolan
Ir. Michel Bagein



Année académique 2017 - 2018
Version 2.0



Informatique Parallèle et Distribuée

Notes de laboratoire

Pr. Manneback Pierre
Dr. Sidi Ahmed Mahmoudi
Ir. François Rolan
Ir. Michel Bagein



Année académique 2017 - 2018
Version 2.0

Contents

1	OpenMP	3
1.1	Introduction	3
1.2	Présentation	3
1.3	Utilité	3
1.4	Les directives OpenMP	4
1.4.1	Parallel	4
1.4.2	For	5
1.4.3	Sections/Section	7
1.4.4	Task	8
1.4.5	Single	9
1.4.6	Barrier	11
1.4.7	Master	12
1.4.8	Critical	13
1.4.9	Atomic	13
1.4.10	Taskwait	14
1.5	Les clauses OpenMP	14
1.5.1	If	15
1.5.2	Shared	15
1.5.3	Private	15
1.5.4	Firstprivate	16
1.5.5	Lastprivate	17
1.5.6	Reduction	17
1.5.7	Default	18
1.5.8	Schedule	18
1.5.9	Ordered	20
1.5.10	Clause Nowait	20
1.6	Les routines OpenMP	21
2	OpenMPI	23
2.1	Introduction	23
2.1.1	Présentation	23
2.1.2	Communicateurs	24
2.1.3	Compilation et exécution	24
2.1.4	Appel bloquant - Appel non bloquant	25
2.2	Librairie de fonctions	25
2.2.1	MPI_Init	25
2.2.2	MPI_Finalize	26
2.2.3	MPI_Get_processor_name	27

2.2.4	MPI_Comm_size	27
2.2.5	MPI_Comm_rank	28
2.2.6	MPI_Send	29
2.2.7	MPI_Isend	30
2.2.8	MPI_Recv	32
2.2.9	MPI_Irecv	34
2.2.10	MPI_Bcast	36
2.2.11	MPI_Reduce	37
2.2.12	MPI_Scatter	39
2.2.13	MPI_Gather	41
2.2.14	MPI_Wait	41
2.2.15	MPI_Test	43
2.2.16	MPI_Probe	44
2.2.17	MPI_Iprobe	46
2.2.18	MPI_Get_count	48
2.2.19	MPI_Barrier	49
2.2.20	MPI_Wtime	50
3	Les processeurs graphiques et CUDA	53
3.1	Introduction	53
3.2	Architecture des processeurs centraux et graphiques	53
3.3	Structures des mmoires ddies aux processeurs graphiques	54
3.4	Les langages de programmation GPU	56
3.4.1	OpenGL	56
3.4.2	CUDA : programmation des processeurs graphiques NVIDIA	57
3.4.3	OpenCL	66
3.5	Exploitation des architectures multi-cœurs htrognes	67
3.5.1	StarPU	67
3.5.2	StarSs	68
3.5.3	GrandCentralDispatch	68

Chapter 1

OpenMP

1.1 Introduction

Ce chapitre présente OpenMP et ses différentes fonctionnalités. Après une brève introduction à OpenMP, nous présentons les différents composants qu'offre cette interface de programmation :

1. les directives OpenMP
2. les clauses OpenMP
3. les routines OpenMP

1.2 Présentation

OpenMP [1] est une interface de programmation pour le multithreading. Elle facilite la programmation d'applications parallèles à mémoire partagée en langage C/C++ ou Fortran depuis 1997. Les spécifications de cette interface sont développées par l'association sans but lucratif "OpenMP Architecture Review Board" (OpenMP ARB) [1] et implémentées par la librairie GOMP [2] incluse dans le compilateur GCC.

1.3 Utilité

La parallélisation d'un code source avec OpenMP s'effectue à l'aide de directives de compilation. Une directive de compilation est une ligne du code source qui débute par `#pragma` et qui contient des mots-clés indiquant au compilateur des actions à effectuer lors de la compilation. Toutes les directives propres à OpenMP commencent par le mot-clé `omp` (`#pragma omp`) et portent des indications sur la manière dont l'application doit être parallélisée. Cette parallélisation est effectuée par défaut à l'aide de threads POSIX (paradigme de mémoire partagée).

La parallélisation peut être spécifiée de trois manières :

1. Les boucles : la directive `#pragma omp for` spécifie que les itérations d'une boucle `for` doivent être exécutées en parallèle par plusieurs threads.
2. Les sections : les directives `#pragma omp sections` et `#pragma omp section` spécifient un ensemble de sections qui doivent être exécutées en parallèle.
3. Les tâches : la directive `#pragma omp task` délimite une partie du code qui doit être encapsulée dans une tâche et exécutée en parallèle par un thread différent du thread appelant.

Des clauses peuvent être jointes à ces directives pour affiner la parallélisation :

- **Private** : les threads générés se partagent le même espace mémoire et donc les mêmes variables. `private` identifie les variables pour lesquelles les threads doivent avoir une copie locale.
- **If** : indique une condition à la parallélisation (nombre d'itérations de la boucle, taille d'une variable ...).
- **Reduction** : spécifie une opération de réduction (sommer les éléments d'un tableau réparti entre plusieurs threads dans un même emplacement mémoire ...).
- **Schedule** : spécifie l'ordonnancement des itérations de la boucle parallélisée.
- ...

Il existe également des directives influençant les flux d'exécution, comme :

- **Barrier** : place une barrière dans le code à laquelle tous les threads se synchronisent.
- **Master** : indique que seul le thread maître peut exécuter la section délimitée par la directive.
- **Critical** : délimite une section de code qui ne peut être exécutée que par un seul thread.
- ...

Enfin, des routines permettent de spécifier le nombre de threads devant être utilisé dans une section parallèle (par défaut, il y a autant de threads que de cœurs processeur), identifier le numéro du thread appelant, créer et utiliser des mutex ...

Le modèle de programmation d'OpenMP utilise le multithreading à mémoire partagée. Ce partage de la mémoire impose au programme de rester au sein d'un ordinateur unique¹, mais cela évite les transferts de données entre les fils d'exécution. Le partage de la mémoire peut entraîner l'apparition de sections critiques dans lesquelles les threads doivent être synchronisés pour éviter les conflits.

1.4 Les directives OpenMP

Cette section et les suivantes se basent sur les informations contenues dans [3] et [4].

1.4.1 Parallel

Syntaxe

```
#pragma omp parallel [clauses]
```

Les clauses employables sont présentes la figure 1.1.

Description

La directive permet de spécifier une section de code qui sera exécutée en parallèle par plusieurs threads. Lorsque le processus rencontre cette directive, il crée un ensemble de threads contenant autant de threads que de cœurs processeur. Le thread maître est le thread numéro 0. Les threads se synchronisent à la fin de la section parallèle (par exemple, après la fonction `sum` du listing 1.1).

Les directives comme `for`, `section` ... qui indiquent la manière de paralléliser les calculs doivent toutes se trouver dans un bloc `parallel`.

Il est interdit de sortir du bloc `parallel` à l'aide d'une instruction `break` ou autre.

¹Sauf mise en place de mécanisme particulier comme un espace de mémoire global ou virtuellement partagé.

Illustration

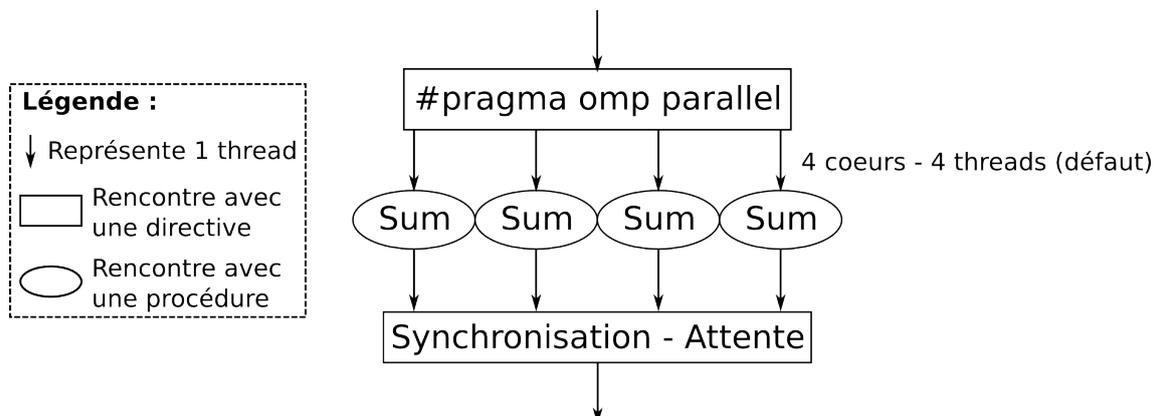


Figure 1.1: Parallel - Quatre threads effectuent chacun la fonction `sum`

Listing 1.1: Illustration de la directive `parallel` : somme des éléments d'un tableau

```

1 void sum(){
2     int i, sum = 0;
3     for (i=0; i<1000; ++i)
4         sum += i;
5     printf("Thread %d a calcul la somme %d\n ", omp_get_thread_num(), sum);
6 }
7
8 int main(){
9     #pragma omp parallel
10    sum();
11    printf("Fin\n");
12 }
13
14 [fremals@fermi01 OpenMP]$ ./parallel_exe
15 Thread 0 a calcul la somme 499500
16 Thread 3 a calcul la somme 499500
17 Thread 2 a calcul la somme 499500
18 Thread 1 a calcul la somme 499500
19 Fin

```

1.4.2 For

Syntaxe

```
#pragma omp for [clauses]
```

Les clauses employables sont présentes la figure 1.1.

Description

La directive `for` spécifie que les itérations d'une boucle doivent être exécutées en parallèle par l'ensemble de threads. Par défaut, OpenMP divise l'ensemble des itérations en autant d'ensembles consécutifs qu'il y a de threads. Chaque ensemble est exécuté par un thread. Une barrière implicite se situe à la fin du bloc `for`. La valeur du gardien de boucle ne peut pas changer au cours de l'exécution des itérations de la boucle.

Il est possible d'écrire directement `#pragma omp parallel for` comme présent au listing 1.2. Il est interdit de sortir du bloc `for` l'aide d'une instruction `break` ou autre.

Illustration

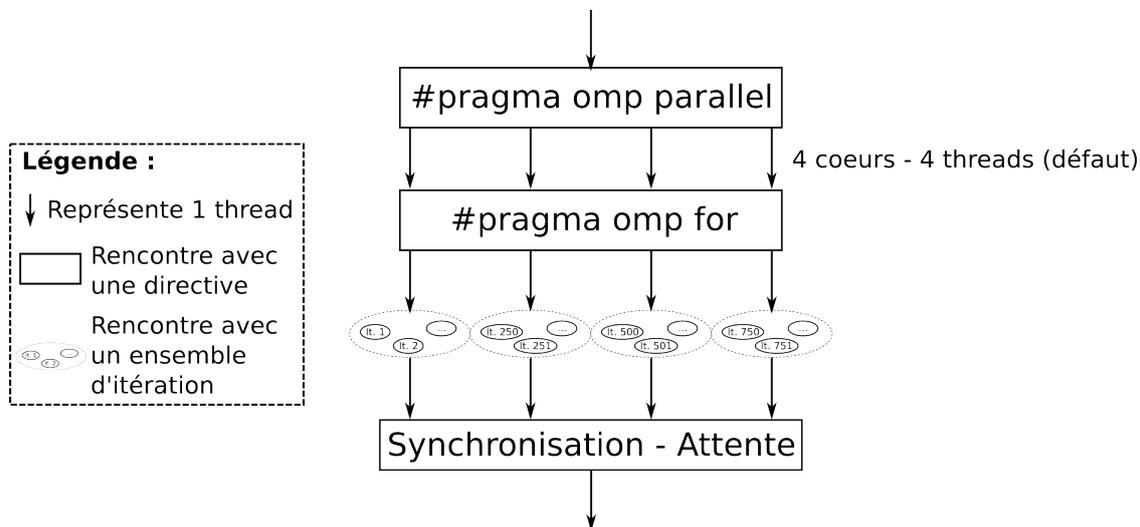


Figure 1.2: For - Quatre threads effectuent chacun 250 itérations de la boucle

Listing 1.2: Illustration de la directive `for` : addition de matrice

```

1 int main(){
2     int i, sum=0;
3     #pragma omp parallel for reduction (+:sum)
4     for (i =0; i<8; i++){
5         sum += i;
6         printf("Thread %d additionne %d\n " , omp_get_thread_num(), i);
7     }
8     printf("Somme finale : %d \n ", sum);
9 }
10
11 Alternative :
12
13 int main(){
14     int i, sum=0;
15     #pragma omp parallel
16     #pragma omp for reduction (+:sum)
17     for (i =0; i<8; i++){
18         sum += i;
19         printf("Thread %d additionne %d\n " , omp_get_thread_num(), i);
20     }
21     printf("Somme finale : %d \n ", sum);
22 }
23
24 [fremals@fermi01 OpenMP]$ ./for_exe
25 Thread 0 additionne 0
26 Thread 0 additionne 1
27 Thread 3 additionne 6
28 Thread 3 additionne 7

```

```

29 Thread 1 additionne 2
30 Thread 1 additionne 3
31 Thread 2 additionne 4
32 Thread 2 additionne 5
33 Somme finale : 28

```

Comme la variable `sum` est partagée, il faut indiquer OpenMP qu'elle fait l'objet d'une opération de réduction. La clause `reduction(+:sum)` indique que chaque thread exécute sa somme locale sur une copie locale de la variable `sum` et, une fois les itérations exécutées, que les sommes locales sont additionnées dans la variable `sum` du thread maître. Pour plus d'information sur cette clause, veuillez consulter la section [1.5.6](#).

1.4.3 Sections/Section

Syntaxe

```
#pragma omp sections [clauses]
```

```
    #pragma omp section
```

Les clauses employables sont présentées à la figure [1.1](#).

Description

Les directives `section` doivent se trouver au sein d'un bloc `sections`. Les différentes `section` appartenant à un même bloc `sections` sont exécutées en parallèle, chacune par un thread. Une barrière implicite se trouve à la fin du bloc `sections`.

Il est interdit de sortir du bloc `section` à l'aide d'une instruction `break` ou autre.

Illustration

Listing 1.3: Illustration des directives `sections` et `section` : addition de matrice

```

1 void computation1(int * tab, int size){
2     int i;
3     for(i=0; i<size; ++i)
4         tab[i]+=i;
5     printf("Thread %d a trait le tableau 1\n ", omp_get_thread_num());
6 }
7
8 void computation2(int * tab, int size){
9     int i;
10    for(i=0; i<size; ++i)
11        tab[i]*=i;
12    printf("Thread %d a trait le tableau 2\n ", omp_get_thread_num());
13 }
14
15 int main(){
16     int tab1[3], tab2[6];
17     int i;
18
19     for(i=0; i<3; ++i){
20         tab1[i] = tab2[i] = i;
21     }
22     for(; i<6; ++i){
23         tab2[i] = i;
24     }
25 }

```

```

26     #pragma omp parallel
27     #pragma omp sections
28     {
29         #pragma omp section
30         computation1(tab1, 3);
31         #pragma omp section
32         computation2(tab2,6);
33     }
34
35     for(i=0;i<3;++i){
36         printf("Tab1[%d] : %d\n", i, tab1[i]);
37     }
38     for(i=0;i<6;++i){
39         printf("Tab2[%d] : %d\n", i, tab2[i]);
40     }
41 }
42
43 [fremals@fermi01 OpenMP]$ ./section_exe
44 Thread 0 a trait le tableau 2
45 Thread 2 a trait le tableau 1
46 Tab1[0] : 0
47 Tab1[1] : 2
48 Tab1[2] : 4
49 Tab2[0] : 0
50 Tab2[1] : 1
51 Tab2[2] : 4
52 Tab2[3] : 9
53 Tab2[4] : 16
54 Tab2[5] : 25

```

1.4.4 Task

Syntaxe

```
#pragma omp task [clauses]
```

Les clauses employables sont présentes la figure 1.1.

Description

La directive `task` limite une partie du code devant être encapsulé dans une tâche. Tous les threads atteignant cette directive créent une nouvelle tâche. Il faut donc veiller à contrôler le nombre de threads atteignant cette directive pour éviter de créer plus de tâches que nécessaire. Les tâches sont exécutées en parallèle par l'ensemble de threads. Une tâche n'est pas forcément exécutée par le thread qui l'a créée.

Il est interdit de sortir du bloc `task` à l'aide d'une instruction `break` ou autre.

Illustration

Listing 1.4: Illustration de la directive `task` : exécution de calculs intensifs en parallèle

```

1 void computation1(int * tab, int size){
2     int i;
3     for(i=0; i<size; ++i)
4         tab[i]+=i;
5     printf("Thread %d a trait le tableau 1\n ", omp_get_thread_num());
6 }
7

```

```

8 void computation2(int * tab, int size){
9     int i;
10    for(i=0; i<size; ++i)
11        tab[i]*=i;
12    printf("Thread %d a trait le tableau 2\n ", omp_get_thread_num());
13 }
14
15 int main(){
16     int tab1[3], tab2[6];
17     int i;
18
19     for(i=0;i<3;++i){
20         tab1[i] = tab2[i] = i;
21     }
22     for(;i<6;++i){
23         tab2[i] = i;
24     }
25
26     omp_set_num_threads(2);
27
28     #pragma omp parallel
29     {
30         #pragma omp task
31         computation1(tab1, 3);
32         #pragma omp task
33         computation2(tab2,6);
34     }
35
36     for(i=0;i<3;++i){
37         printf("Tab1[%d] : %d\n", i, tab1[i]);
38     }
39     for(i=0;i<6;++i){
40         printf("Tab2[%d] : %d\n", i, tab2[i]);
41     }
42 }
43
44 [fremals@fermi01 OpenMP]$ ./task_exe
45 Thread 0 a trait le tableau 1
46 Thread 0 a trait le tableau 1
47 Thread 0 a trait le tableau 2
48 Thread 1 a trait le tableau 2
49 Tab1[0] : 0
50 Tab1[1] : 3
51 Tab1[2] : 6
52 Tab2[0] : 0
53 Tab2[1] : 1
54 Tab2[2] : 8
55 Tab2[3] : 27
56 Tab2[4] : 64
57 Tab2[5] : 125

```

Grce la fonction `omp_set_num_threads`, les sections parallle sont excutes par un ensemble de 2 threads. Chacun des deux threads rencontre les directives `task` et cre deux tches. Nous vrfions que les tches ne sont pas forcement consommées par le thread qui les a cr.

1.4.5 Single

Syntaxe

```
#pragma omp single [clauses]
```

Les clauses employables sont présentes la figure 1.1.

Description

La directive `single` limite une région du code qui est exécutée une seule fois par un seul thread. Une barrière implicite est placée à la fin de cette région.

Il est interdit de sortir du bloc `single` à l'aide d'une instruction `break` ou autre.

Illustration

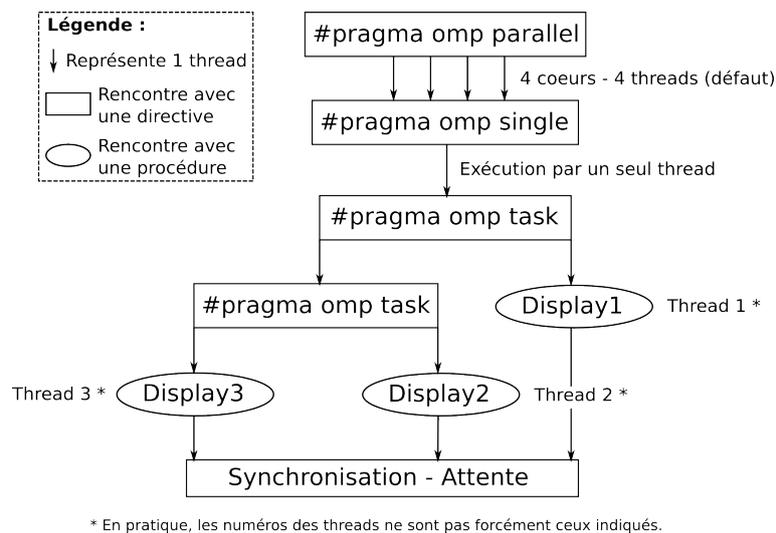


Figure 1.3: Droulement de la parallélisation lors de l'utilisation des directives Task et Single

Listing 1.5: Illustration de la directive `single` : création de tâches sans doublon

```

1 void computation1(int * tab, int size){
2     int i;
3     for(i=0; i<size; ++i)
4         tab[i]+=i;
5     printf("Thread %d a trait le tableau 1\n ", omp_get_thread_num());
6 }
7
8 void computation2(int * tab, int size){
9     int i;
10    for(i=0; i<size; ++i)
11        tab[i]*=i;
12    printf("Thread %d a trait le tableau 2\n ", omp_get_thread_num());
13 }
14
15 int main(){
16     int tab1[3], tab2[6];
17     int i;
18
19     for(i=0; i<3; ++i){
20         tab1[i] = tab2[i] = i;
21     }
22     for(; i<6; ++i){
23         tab2[i] = i;
24     }
25
26     omp_set_num_threads(2);
27

```

```

28     #pragma omp parallel
29     #pragma omp single
30     {
31         #pragma omp task
32         computation1(tab1, 3);
33         #pragma omp task
34         computation2(tab2,6);
35     }
36
37     for(i=0;i<3;++i){
38         printf("Tab1[%d] : %d\n", i, tab1[i]);
39     }
40     for(i=0;i<6;++i){
41         printf("Tab2[%d] : %d\n", i, tab2[i]);
42     }
43 }
44
45 [fremals@fermi01 OpenMP]$ ./single_exe
46 Thread 1 a trait le tableau 1
47 Thread 0 a trait le tableau 2
48 Tab1[0] : 0
49 Tab1[1] : 2
50 Tab1[2] : 4
51 Tab2[0] : 0
52 Tab2[1] : 1
53 Tab2[2] : 4
54 Tab2[3] : 9
55 Tab2[4] : 16
56 Tab2[5] : 25

```

La section crant les tches est excute par un seul thread et une seule tche est cre pour chacun des calculs. Les tches sont ensuite rptaries entre les threads existant.

1.4.6 Barrier

Syntaxe

```
#pragma omp barrier
```

Description

La directive **barrier** synchronise les threads : lorsqu'un thread rencontre cette directive, il arrte son excution jusqu' ce que tous les threads l'aient atteint. Tous les threads de la section parallle doivent donc rencontrer cette directive sous peine de blocage du programme.

Illustration

Listing 1.6: Illustration de la directive **barrier** : mise en vidence des barrire implicites

```

1 int main(){
2     int i, sum=0;
3     #pragma omp parallel
4     {
5         double startTime = omp_get_wtime();
6         while( (omp_get_wtime() - startTime) < (double)(omp_get_thread_num()));
7         printf("%f - Thread %d c o m p l t \n", omp_get_wtime(), omp_get_thread_num());
8         #pragma omp barrier
9         printf("%f - B a r r i e r e   p a s s e \n", omp_get_wtime());
10    }

```

```
11 }
12
13 [fremals@fermi01 OpenMP]$ ./barrier_exe
14 15599.287649 - Thread 0 c o m p l t
15 15600.287649 - Thread 1 c o m p l t
16 15601.287648 - Thread 2 c o m p l t
17 15602.287649 - Thread 3 c o m p l t
18 15602.287680 - B a r r i e p a s s e
19 15602.287686 - B a r r i e p a s s e
20 15602.287681 - B a r r i e p a s s e
21 15602.287757 - B a r r i e p a s s e
```

1.4.7 Master

Syntaxe

```
#pragma omp master
```

Description

La directive `master` delimitte une région qui ne sera exécutée que par le thread maître (les autres threads sautant la région et continuant exécuter le reste du code). Il n'y a pas de barrière implicite à la fin de la région.

Il est interdit de sortir du bloc `master` l'aide d'une instruction `break` ou autre.

Illustration

Listing 1.7: Illustration de la directive `Master` : affichage de résultats

```
1 int main(){
2     int i, a[5];
3     #pragma omp parallel
4     {
5         #pragma omp for
6         for (i = 0; i < 5; i++)
7             a[i] = i * i;
8         #pragma omp master
9         for (i = 0; i < 5; i++)
10            printf("a[%d] = %d\n", i, a[i]);
11        #pragma omp barrier
12        #pragma omp for
13        for (i = 0; i < 5; i++)
14            a[i] += i;
15    }
16 }
17
18 [fremals@fermi01 OpenMP]$ ./master_exe
19 a[0] = 0
20 a[1] = 1
21 a[2] = 4
22 a[3] = 9
23 a[4] = 16
```

1.4.8 Critical

Syntaxe

```
#pragma omp critical (name)
```

Description

La directive `critical` limite une région de code qui sera exécutée par un seul thread à la fois. Lorsqu'un thread atteint une section critique en cours d'exécution par un autre thread, il attend jusqu'à ce que ce dernier soit sorti de la région.

Les sections critiques peuvent être nommées. Plusieurs sections critiques nommées de la même manière sont traitées comme une seule et même section. Les sections critiques portant des noms différents peuvent être exécutées en parallèle. Toutes les sections critiques non nommées sont traitées comme une même section.

Il est interdit de sortir du bloc `critical` à l'aide d'une instruction `break` ou autre.

Illustration

Listing 1.8: Illustration de la directive `Critical` : incrémentation d'une variable

```
1 int main(){
2     int x = 1, y=2;
3     #pragma omp parallel
4     {
5         #pragma omp critical (incr)
6         x += x;
7
8         #pragma omp critical (mul)
9         y *= y;
10    }
11    printf("x : %d, y : %d\n", x, y);
12 }
13
14 [fremals@fermi01 OpenMP]$ ./critical_exe
15 x : 16, y : 65536
```

1.4.9 Atomic

Syntaxe

```
#pragma omp atomic
```

Description

La directive `atomic` permet de spécifier la présence d'une opération demandant une mise à jour atomique de la mémoire. Cette directive fournit ainsi une section critique plus restreinte que celle procurée par `critical`.

L'opération de mise à jour ne peut être quelconque, elle doit obéir à ces règles :

1. au niveau syntaxique, elle s'écrit `var op= expr, opop var` ou `var opop`
2. l'opérateur `op` peut être :
`+, -, *, /, ^, \&, |, <<, >>`
3. la variable doit être scalaire

Illustration

Listing 1.9: Illustration de la directive `atomic` : incréments d'une variable

```

1 #include <stdio.h>
2 #include <omp.h>
3
4 int main(){
5     int x = 1;
6     #pragma omp parallel
7     {
8         #pragma omp atomic
9         ++x;
10    }
11    printf("x : %d\n", x);
12 }
13
14 [fremals@fermi01 OpenMP]$ ./atomic_exe
15 x : 5

```

1.4.10 Taskwait

Syntaxe

```
#pragma omp taskwait
```

Description

La directive `taskwait` spécifie que la tâche qui la rencontre attend l'achèvement des tâches filles gènes depuis le début de la tâche courante.

1.5 Les clauses OpenMP

Les clauses OpenMP permettent d'affiner la parallélisation. Leurs compatibilités avec les directives sont présentées à la figure 1.1.

Table 1.1: Liste incomplète des clauses pouvant être appliquées aux directives

Clauses	Parallel	For	Sections	Single	Parallel For	Task
If	•				•	•
Private	•	•	•	•	•	•
Shared	•	•			•	•
Firstprivate	•	•	•	•	•	•
Lastprivate		•	•		•	
Reduction	•	•	•		•	
Default	•				•	•
Schedule		•			•	
Ordered		•			•	
Nowait		•	•	•		
Num_threads	•	•	•		•	

1.5.1 If

Syntaxe

```
if(condition)
```

Description

La clause `if` conditionne la parallisation : si l'expression de la clause est vraie, la région associée à la directive est exécutée séquentiellement par un thread unique.

Illustration

Listing 1.10: Illustration de la clause `if`

```
1 int main(){
2     int cond = 0;
3     while(cond < 2){
4         #pragma omp parallel if (cond)
5         if (omp_in_parallel()){
6             #pragma omp single
7             printf("Condition : %d, parallelis (%d threads).\n", cond,
8                 omp_get_num_threads());
9         }
10        else{
11            printf("Condition : %d, serialis (%d thread).\n", cond,
12                omp_get_num_threads());
13        }
14        ++cond;
15    }
16 }
17
18 [fremals@fermi01 OpenMP]$ ./if_exe
19 Condition : 0, serialis (1 thread).
20 Condition : 1, parallelis (4 threads).
```

1.5.2 Shared

Syntaxe

```
shared (liste)
```

Description

La clause `shared(liste)` indique la `liste` des variables qui doivent être partagées entre les différents threads. Les variables sont par défaut partagées. Une variable partagée se trouve dans un seul espace mémoire qui est accessible par tous les threads. Les variables partagées peuvent engendrer des sections critiques dont l'accès doit être géré par le programmeur.

1.5.3 Private

Syntaxe

```
private (liste)
```

Description

La clause `private(liste)` indique la liste des variables qu'il faut rendre privées. Lorsqu'une variable est privée, chaque thread en possède une copie locale. C'est-à-dire qu'une nouvelle variable, non initialisée et du même type que la variable privée, est utilisée au sein de chaque thread à la place de la variable privée. Toutes les références à la variable originale sont remplacées par les références à la nouvelle variable. La valeur d'une variable privée n'est pas accessible en dehors de la zone dans laquelle elle est déclarée.

Illustration

Listing 1.11: Illustration de la clause `private` : affichage du numéro des threads

```

1 int main(){
2     int pr1 = 921, pr2 = 473;
3     #pragma omp parallel private(pr1, pr2)
4     {
5         #pragma omp master
6         printf("pr1 : %d, pr2 : %d.\n", pr1, pr2);
7         #pragma omp barrier
8         pr1 = omp_get_thread_num();
9         pr2 = -omp_get_thread_num();
10        printf("%d - pr1 : %d, pr2 : %d\n", omp_get_thread_num(), pr1, pr2);
11    }
12    printf("pr1 : %d, pr2 : %d\n", pr1, pr2);
13 }
14
15 [fremals@fermi01 OpenMP]$ ./private_exe
16 pr1 : 32767, pr2 : 1535201808.
17 0 - pr1 : 0, pr2 : 0
18 3 - pr1 : 3, pr2 : -3
19 1 - pr1 : 1, pr2 : -1
20 2 - pr1 : 2, pr2 : -2
21 pr1 : 921, pr2 : 473

```

1.5.4 Firstprivate

Syntaxe

`firstprivate (list)`

Description

La clause `firstprivate` rend privées les variables de la `list` et les initialise à la valeur qu'elles possèdent avant d'entrer dans la région concernée.

Illustration

Listing 1.12: Illustration de la clause `firstprivate` : affichage du numéro des threads

```

1 int main(){
2     int pr1 = 921, pr2 = 473;
3     #pragma omp parallel firstprivate(pr1, pr2)
4     {
5         #pragma omp master
6         printf("pr1 : %d, pr2 : %d.\n", pr1, pr2);

```

```

7     }
8 }
9
10 [fremals@fermi01 OpenMP]$ ./firstprivate_exe
11 pr1 : 921, pr2 : 473.
12 }

```

1.5.5 Lastprivate

Syntaxe

lastprivate (list)

Description

La clause `lastprivate` rend privées les variables de la `list` et copie la valeur que possède la variable privée la dernière itération ou section dans la variable d'origine.

Illustration

Listing 1.13: Illustration de la clause `lastprivate` : affichage du numéro des threads

```

1 int main(){
2     const int N = 10;
3     double x, a[N];
4     int i;
5
6     #pragma omp parallel for lastprivate(x)
7     for(i=0; i<=N; ++i){
8         x = sin(i);
9         a[i] = exp(x);
10        printf(" I t ration  %d : %f - %f\n", i, x, a[i]);
11    }
12    printf("x : %f\n", x);
13 }
14
15 [fremals@fermi01 OpenMP]$ ./lastprivate_exe
16 I t ration  0 : 0.000000 - 1.000000
17 I t ration  1 : 0.841471 - 2.319777
18 I t ration  4 : -0.756802 - 0.469164
19 I t ration  5 : -0.958924 - 0.383305
20 I t ration  2 : 0.909297 - 2.482578
21 I t ration  3 : 0.141120 - 1.151563
22 x : -0.958924
23 }

```

1.5.6 Reduction

Syntaxe

reduction (opérateur : liste)

Description

La clause `reduction` spécifie que les variables de la liste sont privées et qu'une opération de réduction, spécifiée par l'opérateur, aura lieu sur elles à la fin du bloc.

Les opérateurs possibles pour la réduction sont :

+, *, -, &, ^, |, &&, or ||

Illustration

Listing 1.14: Illustration de la clause `reduction`

```
1 int main(){
2     int sum, i, N=10;
3
4     #pragma omp parallel for reduction(+:sum)
5     for (i=0; i<N; i++) {
6         sum += i;
7     }
8     printf("sum : %d\n", sum);
9 }
10
11 [fremals@fermi01 OpenMP]$ ./reduction_exe
12 sum : 45
```

1.5.7 Default

Syntaxe

`default(shared|none)`

Description

La clause `default` spécifie la portée par défaut des variables au sein de la région parallèle. La clause `default(shared)` garde les variables partagées et la clause `default(none)` entraîne une erreur de compilation si une variable est utilisée dans la région parallèle sans que sa portée n'ait été définie à l'aide d'une des clauses `Private`, `Shared`, `Firstprivate`, `Lastprivate` ou `Reduction`.

1.5.8 Schedule

Syntaxe

`schedule(type [, taille])`

Description

La clause `schedule` spécifie l'ordonnement de l'exécution des itérations d'une boucle parallèle. Il y a 4 types d'ordonnement possibles :

1. **static** : l'ensemble des itérations est divisé en blocs de taille spécifiée par `taille` (si rien n'est spécifié, les itérations sont divisées en autant de blocs qu'il y a de threads pour les traiter ; les tailles de ces blocs peuvent être légèrement différentes) et les blocs sont assignés statiquement aux différents threads selon l'algorithme Round-Robin dans l'ordre des numéros des threads.
2. **dynamic** : l'ensemble des itérations est divisé en blocs de taille spécifiée par `taille` (si rien n'est spécifié, la taille par défaut est 1) et les blocs sont dynamiquement attribués aux threads en attente (lorsqu'un thread a traité un bloc, il attend d'être assigné un autre bloc jusqu'à ce qu'il n'y en ait plus).

3. `guided` : la taille des blocs d'itération est dynamique ; si la `taille` est mise à 1, la taille des blocs d'itération est proportionnelle au nombre d'itérations non assignées divisé par le nombre de threads. Si la `taille` est mise à la valeur k , la taille des blocs est de minimum k itérations, sauf pour le bloc final. La `taille` par défaut est 1.
4. `runtime` : l'ordonnancement est décidé lors de l'exécution via la variable d'environnement `OMP_SCHEDULE`. Il n'est pas possible de spécifier une `taille`.

Illustration

Listing 1.15: Illustration de la clause `schedule`

```

1 #define NUM_LOOPS 18
2 #define SLEEP_EVERY_N 6
3 #define STATIC_CHUNK 3
4
5 #define DYNAMIC_CHUNK 3
6
7 int main(int argc, char * argv){
8     int i, nStatic[NUM_LOOPS], nDynamic[NUM_LOOPS], nGuided[NUM_LOOPS];
9
10    #pragma omp parallel
11    {
12        #pragma omp for schedule(static, STATIC_CHUNK)
13        for (i=0; i<NUM_LOOPS; ++i)
14            nStatic[i] = omp_get_thread_num( );
15
16        #pragma omp for schedule(dynamic, DYNAMIC_CHUNK)
17        for (i=0; i<NUM_LOOPS; ++i)
18            nDynamic[i] = omp_get_thread_num( );
19
20        #pragma omp for schedule(guided)
21        for (i=0; i<NUM_LOOPS; ++i)
22            nGuided[i] = omp_get_thread_num( );
23    }
24
25    printf("Static : ");
26    for(i=0; i<NUM_LOOPS; ++i)
27        printf("%d ", nStatic[i]);
28    printf("\n");
29
30    printf("Dynamic : ");
31    for(i=0; i<NUM_LOOPS; ++i)
32        printf("%d ", nDynamic[i]);
33    printf("\n");
34
35    printf("Guided : ");
36    for(i=0; i<NUM_LOOPS; ++i)
37        printf("%d ", nGuided[i]);
38    printf("\n");
39 }
40
41 [fremals@fermi01 OpenMP]$ ./schedule_exe
42 Static : 0 0 0 1 1 1 2 2 2 3 3 3 0 0 0 1 1 1
43 Dynamic : 0 0 0 2 2 2 1 1 1 3 3 3 0 0 0 3 3 3
44 Guided : 1 1 1 1 1 0 0 0 0 3 3 3 1 1 1 0 0 1
45 }

```

1.5.9 Ordered

Syntaxe

```
#pragma omp ordered
```

Description

La clause / directive `ordered` permet de spécifier que les itérations, ou une partie de celles-ci, doivent être exécutées séquentiellement. Avant d'exécuter ses itérations, un thread doit attendre que les itérations précédentes soient achevées.

Illustration

Listing 1.16: Illustration de la clause `ordered`

```
1 int main(){
2     const int N = 5;
3     double x, a[N];
4     int i;
5
6     #pragma omp parallel for lastprivate(x) ordered
7     for(i=0; i<=N; ++i){
8         x = sin(i);
9         a[i] = exp(x);
10        printf(" I t ration  %d : %f - %f\n", i, x, a[i]);
11    }
12
13    printf("x : %f\n", x);
14 }
15
16 [fremals@fermi01 OpenMP]$ ./ordered_exe
17 I t ration  0 : 0.000000 - 1.000000
18 I t ration  1 : 0.841471 - 2.319777
19 I t ration  2 : 0.909297 - 2.482578
20 I t ration  3 : 0.141120 - 1.151563
21 I t ration  4 : -0.756802 - 0.469164
22 I t ration  5 : -0.958924 - 0.383305
23 x : -0.958924
```

1.5.10 Clause Nowait

Syntaxe

```
nowait
```

Description

La clause `nowait` permet d'annuler la barrière implicite en fin de bloc. Une fois qu'un thread n'a plus de travail, il ne se synchronise pas avec les autres threads et continue son exécution.

Listing 1.17: Illustration de la clause `nowait`

```
1 int main(){
2     int i = 0, N = 3;
3
4     #pragma omp parallel
```

```
5     {
6         #pragma omp for
7         for (i=0; i<N; ++i)
8             printf("Boucle 1, i t ration %d\n", i);
9
10        #pragma omp for nowait
11        for (i=0; i<N; ++i)
12            printf("Boucle 2, i t ration %d\n", i);
13
14        #pragma omp for
15        for (i=0; i<N; ++i)
16            printf("Boucle 3, i t ration %d\n", i);
17    }
18 }
19
20 [fremals@fermi01 OpenMP]$ ./nowait_exe
21 Boucle 1, i t ration 0
22 Boucle 1, i t ration 2
23 Boucle 1, i t ration 1
24 Boucle 2, i t ration 0
25 Boucle 3, i t ration 0
26 Boucle 2, i t ration 2
27 Boucle 3, i t ration 2
28 Boucle 2, i t ration 1
29 Boucle 3, i t ration 1
```

1.6 Les routines OpenMP

De nombreuses routines permettent de paramétrer et d'obtenir des informations sur la parallélisation durant l'exécution du programme (le header `omp.h` doit être inclus pour pouvoir les utiliser) :

- `void omp_set_num_threads(int num_threads)` : spécifie la taille de l'ensemble de threads (entier positif) utilisé pour exécuter les régions de code parallèles. Cette routine doit être appelée depuis une section de code séquentiel.
- `int omp_get_num_threads(void)` : renvoie le nombre de threads en cours d'exécution.
- `int omp_get_thread_num(void)` : renvoie le numéro du thread appelant (le thread maître a le numéro 0).
- `int omp_get_num_procs(void)` : renvoie le nombre de coeurs, d'unités de travail que possède la machine.
- `void omp_set_dynamic(int dynamic_threads)` : permet d'activer ou de désactiver la possibilité d'ajuster la taille de l'ensemble de threads exécutant les régions parallèles (le mode dynamique est activé si la variable `dynamic_threads` est positive non nulle). Cette routine doit être appelée depuis une section de code séquentiel.
- `int omp_get_dynamic(void)` : renvoie 0 si l'ajustement dynamique des threads est activé et une valeur différente de zéro sinon.
- `double omp_get_wtime(void)` : renvoie le nombre de secondes écoulées depuis un certain point dans le passé
- `int omp_in_parallel(void)` : indique si la section de code exécutée est parallélisée (la valeur de retour est différente de zéro si c'est le cas et égale à zéro sinon).

- `void omp_init_lock(omp_lock_t *lock)` : initialise un verrou associé à la variable `lock` 1.18.
- `void omp_destroy_lock(omp_lock_t *lock)` : désassocie la variable `lock` de tout verrouillage.
- `void omp_set_lock(omp_lock_t *lock)` : force le thread appelant à attendre jusqu'à ce que le verrou soit libre (la variable `lock` doit avoir été au préalable initialisé).
- `void omp_unset_lock(omp_lock_t *lock)` : libère le verrou du thread appelant.
- `int omp_test_lock(omp_lock_t *lock)` : tente de verrouiller le verrou mais ne bloque pas si celui-ci l'est déjà.

Listing 1.18: Illustration du verrouillage au sein d'OpenMP

```
1 int main() {
2     omp_lock_t my_lock;
3     omp_init_lock(&my_lock);
4
5     #pragma omp parallel num_threads(4)
6     {
7         int tid = omp_get_thread_num( );
8         int i, j;
9
10
11         omp_set_lock(&my_lock);
12         printf("Thread %d - starting locked region\n", tid);
13         sleep(1);
14         printf("Thread %d - ending locked region\n", tid);
15         omp_unset_lock(&my_lock);
16
17     }
18
19     omp_destroy_lock(&my_lock);
20 }
21
22 [fremals@fermi01 OpenMP]$ ./lock_exe
23 Thread 2 - Entrant dans la r gion critique
24 Thread 2 - Sortant de la r gion critique
25 Thread 3 - Entrant dans la r gion critique
26 Thread 3 - Sortant de la r gion critique
27 Thread 0 - Entrant dans la r gion critique
28 Thread 0 - Sortant de la r gion critique
29 Thread 1 - Entrant dans la r gion critique
30 Thread 1 - Sortant de la r gion critique
```

Chapter 2

OpenMPI

2.1 Introduction

Ce chapitre introduit l'utilisation d'OpenMPI. Il commence par présenter brièvement MPI et OpenMPI puis présente certaines fonctions OpenMPI.

2.1.1 Présentation

MPI (Message Passing Interface) est une API¹ spécifiant la parallélisation de processus sur des systèmes mémoire distribuée. Les communications entre processus s'effectuent par mission/réception de messages. Le projet a commencé en 1992 et en est aujourd'hui la version 2.0.

Open MPI est un projet open-source, fruit de la fusion de plusieurs projets, qui a vu le jour en 2003. Le but recherché au travers cette fusion est de fournir la meilleure bibliothèque MPI. Ce projet permet d'intégrer les fonctions MPI des codes C/C++ ou Fortran, de compiler les codes modifiés et d'exécuter les programmes sur des systèmes multi-cœurs mémoire distribuée.

Le but de cette API est de permettre la parallélisation sur un ordinateur, sur plusieurs ordinateurs, sur plusieurs nœuds d'un cluster ... Contrairement à OpenMP, la mémoire est distribuée et les multiples processus MPI possèdent leur propre mémoire. Les communications inter-processus sont assurées par le paradigme de passage de messages (send/receive), ceci rendant les transferts de données plus coûteux qu'avec OpenMP. C'est pourquoi il peut parfois être intéressant de mélanger OpenMP et MPI au sein d'un même programme : MPI distribue alors les grosses tâches entre les différentes machines, et OpenMP parallélise ces tâches au sein de chaque machine.

OpenMPI offre des méthodes permettant l'initialisation et la terminaison des processus MPI. La portion de code comprise entre ces deux fonctions sera exécutée autant de fois qu'il y a de processus. L'utilisateur utilise le rang des processus pour filtrer et attribuer les traitements aux processus. Généralement, un mécanisme maître/esclave est utilisé de telle manière que le maître partage les données à traiter entre les esclaves (et lui-même) et intègre finalement le résultat final.

De nombreuses fonctions permettant l'échange de données entre les processus ont été implémentées. Il est ainsi possible d'envoyer des messages d'un metteur vers un destinataire de manière synchrone, de manière asynchrone, d'envoyer des messages d'un metteur vers un ensemble de processus (broadcasting) ou d'effectuer l'inverse (opération de réduction) ... Il y a également une fonction permettant la synchronisation des processus.

¹Application Programming Interface - Interface de programmation



2.1.2 Communicateurs

OpenMPI permet d'organiser des groupes de processus grce des communicateurs. Les communicateurs permettent d'organiser les communications entre les processus MPI. Ils permettent de crer des groupes de processus et ainsi de dfinir les destinataires ou les metteurs de messages pour des fonctions telles que `MPI_Reduce`, `MPI_Bcast` ...

Dans le cadre de ce cours, nous n'tudierons pas la gestion des communicateurs². Voici un ensemble de communicateurs prexistant dans la librairie OpenMPI et qui pourront servir lors de l'emploi des fonctions de communication :

- `MPI_COMM_WORLD` : communicateurs contenant tous les processus MPI.
- `MPI_COMM_SELF` : contient uniquement le processus appelant (communication intra-processus).
- `MPI_COMM_NULL` : communicateur NULL (utile pour les fonctions o il est permis de passer un communicateur NULL).

2.1.3 Compilation et exécution

La compilation d'un programme utilisant OpenMPI est particulire car elle ncessite l'emploi de la commande `mpicc`. Cette commande englobe le compilateur `gcc` et permet de compiler du code C et de lier au programme les librairies MPI ncessaires. Pour compiler un code, il suffit d'utiliser la ligne de commandes :

```
mpicc prog.c -o prog
```

Pour excuter un programme, il faut utiliser la commande `mpirun`. Ce script inspecte le systme et dmarre les jobs MPI en fonction du matriel disponible. La commande :

```
mpirun -np 2 prog
```

permet de lancer un programme `prog` qui sera excut par deux processus (`np` = number of processes).

Il est possible d'excuter des processus sur une autre machine que celle appelant le programme. Pour ce faire, il suffit d'indiquer le nom des machines htes cibles :

```
mpirun -np 2 --host host1,host2 prog
```

```
mpirun -np 2 --hostfile machines prog
```

L'option `--hostfile` requiert l'emploi d'un fichier, ici nomm `machines`, qui contient le nom des htes (un hte par ligne).

Comme les transferts de donnes s'effectue via une connexion SSH, il est intressant d'installer une cl SSH sur les htes. Ceci permet d'viter de devoir saisir le mot de passe ncessaire la connexion chaque excution. Voici la procdure suivre pour placer une cl SSH sur un hte :

```
[xxx@s-iglab-07a ~]$ ssh-keygen % cration de la cl
Generating public/private rsa key pair.
Enter file in which to save the key (/home/xxx/.ssh/id_rsa): % saisir Entre
Created directory '/home/xxx/.ssh'.
```

²Les personnes intresses peuvent examiner la fonction `MPI_Comm_create` qui permet la cration de communicateurs.



```

Enter passphrase (empty for no passphrase): % saisir Entre
Enter same passphrase again: % saisir Entre
Your identification has been saved in /home/xxx/.ssh/id_rsa.
Your public key has been saved in /home/xxx/.ssh/id_rsa.pub.
The key fingerprint is:
8e:e0:51:29:64:4d:27:60:c1:37:18:f0:00:e0:eb:e2 temp@s-iglab-07a
[xxx@s-iglab-07a ~]$ ssh-copy-id xxx@s-iglab-08a % envoie de la cl sur l'hte
s-iglab-08a
The authenticity of host 's-iglab-08a (10.107.45.38)' can't be established.
RSA key fingerprint is 60:53:7e:94:a1:7a:43:ab:e1:48:fa:61:64:00:68:75.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 's-iglab-08a,10.107.45.38' (RSA) to the list of
known hosts.
xxx@s-iglab-08a's password: % saisir le mot de passe de la session hte
Now try logging into the machine, with "ssh 'xxx@s-iglab-08a'", and check in:

~/ssh/authorized_keys

```

to make sure we haven't added extra keys that you weren't expecting.

Il faut galemment distribuer le fichier executable sur la machines hte. Pour ce faire, excutez la commande :

```
[xxx@s-iglab-07a ~]$ scp ./prog xxx@s-iglab-07a:~
```

et le fichier sera alors plac dans le dossier de l'utilisateur xxx sur l'hte s-iglab-07a.

2.1.4 Appel bloquant - Appel non bloquant

Lorsqu'une routine est bloquante, son retour se fait une fois sa fonctionnalit excute. Si la routine est non-bloquante, son retour est immdiat, mme si sa fonctionnalit n'est pas compltement ralis.

2.2 Librairie de fonctions

Cette section se base sur les informations contenues dans [5], [6] et [7].

2.2.1 MPI_Init

Syntaxe C

```
#include <mpi.h>
int MPI_Init(int *argc, char ***argv)
```

Description

Initialise l'environnement d'excution de MPI. L'appel cette fonction doit ncessairement preder l'emploi de toute autre fonction d'OpenMPI. OpenMPI ne pouvant tre initialis qu'une fois, toute autre appel cette fonction rsultera en un chec.

Paramtres d'entre

`argc` : pointeur vers le nombre d'arguments.

`argv` : tableau d'arguments.

Valeur de retour

MPI_SUCCESS : la routine s'est exécutée sans erreur.

MPI_ERR_OTHER : indique la tentative d'un second appel à la routine MPI_Init, ce qui n'est pas permis.

Illustration

Listing 2.1: Illustration de la fonction MPI_Init

```
1 #include <mpi.h>
2
3 int main(int argc, char * argv[]){
4     MPI_Init(&argc, &argv);
5     /*Programme principal*/
6     MPI_Finalize();
7 }
```

2.2.2 MPI_Finalize

Syntaxe C

```
#include <mpi.h>
int MPI_Finalize()
```

Description

Tous les processus doivent appeler cette routine avant de se terminer. Une fois l'appel à cette routine effectué, il n'est plus possible d'appeler d'autres fonctions MPI, même MPI_Init ou MPI_Init_thread³.

Dans le cas des systèmes multithread, il est conseillé que le thread appelant MPI_Finalize soit le thread ayant appelé MPI_Init ou MPI_Init_thread.

Valeur de retour

MPI_SUCCESS : la routine s'est exécutée sans erreur.

Illustration

Listing 2.2: Illustration de la fonction MPI_Finalize

```
1 #include <mpi.h>
2
3 int main(int argc, char * argv[]){
4     MPI_Init(&argc, &argv);
5     /*Programme principal*/
6     MPI_Finalize();
7 }
```

³Les exceptions à cette règle sont les fonctions MPI_Get_version (indique la version d'OpenMPI utilisée), MPI_Initialized (indique si MPI_Init a déjà été appelé), et MPI_Finalized (indique si MPI_Finalize a déjà été appelé).



2.2.3 MPI_Get_processor_name

Syntaxe C

```
#include <mpi.h>
int MPI_Get_processor_name(char *name, int *resultlen)
```

Description

Cette routine permet de récupérer le nom du processeur et donc d'identifier la machine hôte. Le paramètre `name` contiendra le nom du processeur et doit au moins pouvoir accueillir `MPI_MAX_PROCESSOR_NAME` caractères. Le second paramètre contient le nombre de caractères du nom du processeur.

Paramètres de sortie

`name` : un identifiant unique pour le nœud de travail.
`resultlen` : la longueur de la chaîne pointée par `name`.

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

Illustration

Listing 2.3: Illustration de la fonction `MPI_Get_processor_name`

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, len = 0;
7     char name[MPI_MAX_PROCESSOR_NAME];
8     MPI_Init(&argc, &argv);
9     MPI_Get_processor_name(name, &len);
10    printf("Processor : %s (%d caractres).\n", name, len);
11    MPI_Finalize();
12 }
13
14 [fremals@fermi01 OpenMPI]$ mpirun -np 1 get_proc
15 Processor : fermi01 (7 caractres).
```

2.2.4 MPI_Comm_size

Syntaxe C

```
#include <mpi.h>
int MPI_Comm_size(MPI_Comm comm, int *size)
```

Description

Indique le nombre de processus se trouvant dans le groupe associé au communicateur `comm`.

Paramètres d'entrée

`comm` : communicateur.

Paramtres de sortie

`size` : nombre de processus dans le groupe `comm`.

Valeur de retour

`MPI_SUCCESS` : la routine s'est excute sans erreur.

`MPI_ERR_COMM` : le communicateur n'est pas valide. Une erreur frquente est d'utiliser un communicateur `NULL`.

`MPI_ERR_ARG` : un argument n'est pas valide.

Illustration

Listing 2.4: Illustration de la fonction `MPI_Comm_size`

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0;
7     MPI_Init(&argc,&argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD,&size);
10    printf("Il y a %d processus dans le groupe\n", size);
11
12    MPI_Finalize();
13 }
14
15 [fremals@fermi01 OpenMPI]$ mpirun -np 2 ./comm_size
16 Il y a 2 processus dans le groupe
17 Il y a 2 processus dans le groupe
```

2.2.5 MPI_Comm_rank

Syntaxe C

```
#include <mpi.h>
int MPI_Comm_rank(MPI_Comm comm, int *rank)
```

Description

Indique le rang (le numro) du processus appelant au sein du groupe de processus associ au communicateur `comm`. Le rang permet d'identifier un processus MPI et de lui assigner un ensemble d'instructions particulier.

Paramtres d'entre

`comm` : communicateur.

Paramtres de sortie

`rank` : rang du processus appelant au sein du groupe de processus associ `comm`.

Valeur de retour

MPI_SUCCESS : la routine s'est exécutée sans erreur.

MPI_ERR_COMM : le communicateur n'est pas valide. Une erreur fréquente est d'utiliser un communicateur NULL.

Illustration

Listing 2.5: Illustration de la fonction MPI_Comm_rank

```

1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, rank = 0;
7     MPI_Init(&argc, &argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD, &size);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11    printf("[%d] Il y a %d processus dans le groupe.\n", rank, size);
12
13    MPI_Finalize();
14 }
15
16 [fremals@fermi01 OpenMPI]$ mpirun -np 2 ./comm_rank
17 [0] Il y a 2 processus dans le groupe.
18 [1] Il y a 2 processus dans le groupe.

```

2.2.6 MPI_Send

Syntaxe C

```

#include <mpi.h>
int MPI_Send(void *buf, int count, MPI_Datatype datatype, int dest, int tag,
MPI_Comm comm)

```

Description

Effectue l'envoi bloquant d'un message. La routine bloque jusqu'au transfert du message du buffer d'envoi `buf` au buffer de réception (la boîte aux lettres du destinataire).

Paramètres d'entrée

`buf` : l'adresse de l'emplacement mémoire contenant le message à envoyer.

`count` : nombre d'éléments de type `datatype` à envoyer (entier non-négatif).

`datatype` : type des éléments envoyés. Les types existants sont : MPI_CHAR, MPI_SHORT, MPI_INT, MPI_LONG, MPI_FLOAT, MPI_DOUBLE, MPI_LLONG, MPI_LLONG_DOUBLE⁴.

`dest` : rang du processus destinataire dans le communicateur `comm`.

`tag` : une étiquette pour le message choisie par le programmeur, permet de filtrer les messages.

`comm` : communicateur.

⁴Il est possible de créer des types de données dérivés, des types de données construits à partir des types fournis par MPI, à l'aide des routines MPI_Type_contiguous, MPI_Type_vector, MPI_Type_indexed, MPI_Type_struct, MPI_Type_commit, MPI_Type_free ...

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

`MPI_ERR_COMM` : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur `NULL`.

`MPI_ERR_COUNT` : l'argument `count` est incorrect.

`MPI_ERR_TYPE` : l'argument `datatype` est incorrect.

`MPI_ERR_TAG` : l'argument `tag` est incorrect. Les tags doivent être non-négatifs ; la plus grande valeur permise pour ce paramètre est disponible dans l'attribut `MPI_TAG_UB`.

`MPI_ERR_RANK` : rang du destinataire incorrect. Le rang doit être compris entre 0 et la taille du groupe du communicateur décrément d'une unité.

Illustration

Listing 2.6: Illustration de la fonction `MPI_Send`

```

1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, rank = 0;
7     MPI_Init(&argc, &argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD, &size);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11    if(rank == 0){
12        int tosend = 5;
13        MPI_Ssend(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
14        printf("[%d] Message post .\n", rank);
15    }
16    if(rank==1){
17        sleep(1);
18        MPI_Status status;
19        int torecv;
20        printf("[%d] Lecture de la boîte au lettre.\n", rank);
21        MPI_Recv(&torecv, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
22        printf("[%d] Message reçu : %d.\n", rank, torecv);
23    }
24    MPI_Finalize();
25 }
26
27 [fremals@fermi01 OpenMPI]$ mpirun -np 2 send
28 [0] Envoie du message.
29 [0] Message post .
30 [1] Lecture de la boîte au lettre.
31 [1] Message reçu : 5.

```

2.2.7 `MPI_Isend`

Syntaxe C

```

#include <mpi.h>
int MPI_Isend(void *buf, int count, MPI_Datatype datatype, int dest, int tag,
MPI_Comm comm, MPI_Request *request)

```

Description

Commence l'envoi non-bloquant d'un message. L'appel alloue un objet représentant la requête de communication et l'associe la variable `request`. Cette variable peut ensuite être utilisée pour connaître l'état de la communication (`MPI_Test`) ou attendre son achèvement (`MPI_Wait`).

L'appel `MPI_Isend` indique au système qu'il peut commencer la copie des données en dehors du buffer d'envoi. Le buffer d'envoi ne doit pas être accédé tant que cette copie n'est pas finie.

Paramètres d'entrée

`buf` : l'adresse de l'emplacement mémoire contenant le message à envoyer.

`count` : nombre d'éléments de type `datatype` à envoyer (entier non-négatif).

`datatype` : type des éléments envoyés. Les types existants sont : `MPI_CHAR`, `MPI_SHORT`, `MPI_INT`, `MPI_LONG`, `MPI_FLOAT`, `MPI_DOUBLE`, `MPI_LONGLONG`, `MPI_LONG_DOUBLE`⁵.

`dest` : rang du processus destinataire dans le communicateur `comm`.

`tag` : une étiquette pour le message choisie par le programmeur et permettant de filtrer les messages à l'arrivée. La plus grande valeur possible pour ce paramètre est contenue dans la constante `MPI_TAG_UB`.

`comm` : nom du communicateur.

Paramètres de sortie

`request` : variable permettant de connaître l'état de la communication.

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

`MPI_ERR_COMM` : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur `NULL`.

`MPI_ERR_COUNT` : l'argument `count` est incorrect.

`MPI_ERR_TYPE` : l'argument `datatype` est incorrect.

`MPI_ERR_TAG` : l'argument `tag` est incorrect. Les tags doivent être non-négatifs.

`MPI_ERR_RANK` : rang du destinataire incorrect. Le rang doit être compris entre 0 et la taille du groupe du communicateur décrémenté d'une unité.

Illustration

Listing 2.7: Illustration de la fonction `MPI_Isend`

```

1 #include <stdio.h>
2 #include <mpi.h>
3
4 int computation(){
5     return rand()%10;
6 }
7
8 int main(int argc, char * argv[])
9 {
10     int size = 0, rank = 0;
11     MPI_Init(&argc, &argv);
12

```

⁵Il est possible de créer des types de données dérivés, des types de données construits à partir des types fournis par MPI, à l'aide des routines `MPI_Type_contiguous`, `MPI_Type_vector`, `MPI_Type_indexed`, `MPI_Type_struct`, `MPI_Type_commit`, `MPI_Type_free` ...

```

13 MPI_Comm_size(MPI_COMM_WORLD,&size);
14 MPI_Comm_rank(MPI_COMM_WORLD,&rank);
15 MPI_Request sendState;
16 if(rank == 0){
17     sleep(1);
18     int tosend = computation(), temp;
19     MPI_Isend(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, &sendState);
20     while(tosend!=9){
21         temp = computation(); //le buffer d'envoi ne peut tre modifi
22         MPI_Wait(&sendState, MPI_STATUS_IGNORE);
23         tosend = temp;
24         MPI_Isend(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, &sendState);
25     }
26     MPI_Wait(&sendState, MPI_STATUS_IGNORE);
27 }
28 if(rank==1){
29     MPI_Status status;
30     int torecv=0;
31     while(torecv!=9){
32         MPI_Recv(&torecv, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
33         printf("Valeur calcule : %d\n", torecv);
34     }
35 }
36 MPI_Finalize();
37 }
38
39 [fremals@fermi01 OpenMPI]$ mpirun -np 2 isend
40 Valeur calcule : 3
41 Valeur calcule : 6
42 Valeur calcule : 7
43 Valeur calcule : 9

```

2.2.8 MPI_Recv

Syntaxe C

```

#include <mpi.h>
int MPI_Recv(void *buf, int count, MPI_Datatype datatype, int source, int tag,
MPI_Comm comm, MPI_Status *status)

```

Description

Effectue la rception bloquante d'un message. La routine bloque jusqu' la rception du message. Un appel Recv peut rendre la main avant l'appel au Send correspondant.

Un message peut tre reu uniquement s'il est adress au processus rcepteur, et si les valeurs de la source, de l'tiquette et du communicateur correspondent celles spcifies dans la fonction de rception. Pour la source et l'tiquette, il existe deux valeurs, MPI_ANY_SOURCE et MPI_ANY_TAG, qui spcifient la rception de messages en provenance de n'importe quelle source (appartenant au groupe du communicateur spcifi) et portant n'importe quelle tiquette.

Un processus peut s'envoyer un message soi-mme, mais cela ncessite l'emploi d'une fonction d'envoi asynchrone pour viter tout blocage.

Paramtres d'entre

count : nombre maximum d'lments recevoir de type **datatype** (entier non-ngatif).

datatype : type des lments reus. Les types existant sont : MPI_CHAR, MPI_SHORT, MPI_INT,

MPI_LONG, MPI_FLOAT, MPI_DOUBLE, MPI_LONG, MPI_LONG_DOUBLE⁶.

dest : rang du processus source dans le communicateur **comm**.

tag : une tiquette pour le message choisie par le programmeur et permettant de filtrer les messages. S'il n'y a pas besoin de filtrer les messages, le tag **MPI_ANY_TAG** peut tre utilis. La plus grande valeur possible pour ce paramtre est contenue dans la constante **MPI_TAG_UB**.

comm : communicateur.

Paramtres de sortie

buf : l'adresse de l'emplacement mmoire o le message doit tre stock.

status : structure contenant des informations sur le message reu :

- **status.MPI_TAG** : l'tiquette du message reu.
- **status.MPI_SOURCE** : l'metteur du message reu.

Valeur de retour

MPI_SUCCESS : la routine s'est excute sans erreur.

MPI_ERR_COMM : le communicateur est incorrect. Une erreur frquente est d'utiliser un communicateur **NULL**.

MPI_ERR_COUNT : l'argument **count** est incorrect.

MPI_ERR_TYPE : l'argument **datatype** est incorrect.

MPI_ERR_TAG : l'argument **tag** est incorrect. Les tags doivent tre non-ngatifs ; la plus grande valeur permise pour ce paramtre est disponible dans l'attribut **MPI_TAG_UB**. Dans le cas de la rception, l'tiquette peut valoir **MPI_ANY_TAG**.

MPI_ERR_RANK : rang de la source incorrect. Le rang doit tre compris entre 0 et la taille du groupe du commicateur dcrment d'une unit. Dans le cas de la rception, le rang peut valoir **MPI_ANY_SOURCE**.

Illustration

Listing 2.8: Illustration de la fonction **MPI_Recv**

```

1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, rank = 0;
7     MPI_Init(&argc,&argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD,&size);
10    MPI_Comm_rank(MPI_COMM_WORLD,&rank);
11    if(rank == 0){
12        int tosend = 5;
13        MPI_Ssend(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
14        printf("[%d] Message post .\n", rank);
15    }
16    if(rank==1){
17        sleep(1);
18        MPI_Status status;
19        int torecv;

```

⁶Il est possible de crer des types de donnees drivs, des types de donnees construits partir des types fournis par MPI, l'aide des routines **MPI_Type_contiguous**, **MPI_Type_vector**, **MPI_Type_indexed**, **MPI_Type_struct**, **MPI_Type_commit**, **MPI_Type_free** ...

```

20     printf("[%d] Lecture de la boite au lettre.\n", rank);
21     MPI_Recv(&torecv, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
22     printf("[%d] Message r e u : %d.\n", rank, torecv);
23 }
24 MPI_Finalize();
25 }
26
27 [fremals@fermi01 OpenMPI]$ mpirun -np 2 recv
28 [1] Lecture de la boite au lettre.
29 [0] Envoie du message.
30 [0] Message post .
31 [1] Message r e u : 5 (source : 0,   tiquette   : 5).

```

2.2.9 MPI_Irecv

Syntaxe C

```

#include <mpi.h>
int MPI_Irecv(void *buf, int count, MPI_Datatype datatype, int source, int tag,
MPI_Comm comm, MPI_Request *request)

```

Description

Lance la rception non bloquante d'un message. L'appel alloue un objet representant une requete de communication et l'associe la variable `request`. Cette variable peut tre ultrieurement utilise pour connaitre l'tat de la communication (`MPI_Test`) ou attendre son achvement (`MPI_Wait`).

L'appel `MPI_Irecv` indique au systme d'exploitation qu'il peut commencer crire les donnees dans le buffer de rception. Le buffer de rception ne doit pas tre accd tant que cette copie n'est pas finie.

Paramtres d'entre

`count` : nombre maximum d'lments de type `datatype` recevoir (entier non-ngatif).

`datatype` : type des lments reus. Les types existant sont : `MPI_CHAR`, `MPI_SHORT`, `MPI_INT`, `MPI_LONG`, `MPI_FLOAT`, `MPI_DOUBLE`, `MPI_LONG`, `MPI_LONG_DOUBLE`⁷.

`dest` : rang du processus source dans le communicateur `comm`.

`tag` : une tiquette pour le message choisie par le dveloppeur et permettant de filtrer les messages. S'il n'y a pas besoin de filtrer les messages, le tag `MPI_ANY_TAG` peut tre utilis. La plus grande valeur possible pour ce paramtre est contenue dans la constante `MPI_TAG_UB`.

`comm` : nom du communicateur.

Paramtres de sortie

`buf` : l'adresse de l'emplacement mmoire o le message doit tre stock.

`status` : structure contenant des informations sur le messages reu :

- `status.MPI_TAG` : l'tiquette du message reu.
- `status.MPI_SOURCE` : l'metteur du message reu.

⁷Il est possible de crer des types de donnees drivs, des types de donnees construits partir des types fournis par MPI, l'aide des routines `MPI_Type_contiguous`, `MPI_Type_vector`, `MPI_Type_indexed`, `MPI_Type_struct`, `MPI_Type_commit`, `MPI_Type_free` ...

Valeur de retour

MPI_SUCCESS : la routine s'est exécutée sans erreur.

MPI_ERR_COMM : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur NULL.

MPI_ERR_COUNT : l'argument `count` est incorrect.

MPI_ERR_TYPE : l'argument `datatype` est incorrect.

MPI_ERR_TAG : l'argument `tag` est incorrect. Les tags doivent être non-négatifs.

MPI_ERR_RANK : rang de la source incorrect. Le rang doit être compris entre 0 et la taille du groupe du communicateur dérivé d'une unité. Dans le cas de la réception, le rang peut valoir MPI_ANY_SOURCE.

Illustration

Listing 2.9: Illustration de la fonction MPI_Irecv

```

1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, rank = 0;
7     MPI_Init(&argc, &argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD, &size);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11    if(rank == 0){
12        sleep(1);
13        int tosend = 5;
14        MPI_Send(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
15    }
16    if(rank==1){
17        MPI_Status status;
18        int torecv;
19        MPI_Request recvState;
20        MPI_Irecv(&torecv, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &recvState);
21        int recved = 0;
22        while(!recved){
23            MPI_Test(&recvState, &recved, MPI_STATUS_IGNORE);
24            printf("Etat de la r ception : %d\n", recved);
25            if(!recved)
26                sleep(0);
27        }
28    }
29    MPI_Finalize();
30 }
31
32 [fremals@fermi01 OpenMPI]$ mpirun -np 2 ./irecv
33 Etat de la r ception : 0
34 Etat de la r ception : 0
35 ...
36 Etat de la r ception : 0
37 Etat de la r ception : 0
38 Etat de la r ception : 1

```



2.2.10 MPI_Bcast

Syntaxe C

```
#include <mpi.h>
int MPI_Bcast(void *buffer, int count, MPI_Datatype datatype, int root,
MPI_Comm comm)
```

Description

Cette routine transmet un message du processus source (de rang `root`) vers tous les processus du communicateur, le processus source y compris. Tous les processus du groupe doivent l'appeler en passant les mme valeurs pour les variables `root` et `comm`. La taille des données (`size*sizeof(datatype)`) doit galemment tre identique pour les diffrents processus. La valeur de la variable `buffer` du processus source est copie dans le `buffer` des processus associs au communicateur.

Paramtres d'entre

`buffer` : pointeur vers la zone mmoire qui contient/reoit les données.
`count` : nombre d'lements transférer.
`datatype` : type des lments reus. Les types existant sont : `MPI_CHAR`, `MPI_SHORT`, `MPI_INT`, `MPI_LONG`, `MPI_FLOAT`, `MPI_DOUBLE`, `MPI_LONG`, `MPI_LONG_DOUBLE`⁸.
`root` : rang du processus source.
`comm` : communicateur.

Paramtres de sortie

`buffer` : pointeur vers la zone mmoire qui contient/reoit les données.

Valeur de retour

`MPI_SUCCESS` : la routine s'est excute sans erreur.
`MPI_ERR_COMM` : le communicateur est incorrect. Une erreur frquente est d'utiliser un communicateur `NULL`.
`MPI_ERR_COUNT` : l'argument `count` est incorrect. L'argument ne peut pas tre ngatif.
`MPI_ERR_TYPE` : l'argument `datatype` est incorrect, le type n'est pas connu de MPI.
`MPI_ERR_BUFFER` : le pointeur `buffer` est incorrect. La zone de mmoire pointe doit tre alloue.
`MPI_ERR_ROOT` : le rang du processus source est incorrect.

Illustration

Listing 2.10: Illustration de la fonction `MPI_Bcast`

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int computation(){
5     return rand()%10;
```

⁸Il est possible de crer des types de données drivs, des types de données construits partir des types fournis par MPI, l'aide des routines `MPI_Type_contiguous`, `MPI_Type_vector`, `MPI_Type_indexed`, `MPI_Type_struct`, `MPI_Type_commit`, `MPI_Type_free` ...

```

6 }
7
8 int main(int argc, char * argv[])
9 {
10     int size = 0, rank = 0;
11     MPI_Init(&argc, &argv);
12
13     MPI_Comm_size(MPI_COMM_WORLD, &size);
14     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
15     MPI_Request sendState;
16     if(rank == 0){
17         int tobcast = 4;
18         MPI_Bcast(&tobcast, 1, MPI_INT, 0, MPI_COMM_WORLD);
19         printf("[%d] Valeur diffuse : %d\n", rank, tobcast);
20     }
21     else if(rank != 0){
22         int tobcast = 0;
23         MPI_Bcast(&tobcast, 1, MPI_INT, 0, MPI_COMM_WORLD);
24         printf("[%d] Valeur diffuse : %d\n", rank, tobcast);
25     }
26     MPI_Finalize();
27 }
28
29 [fremals@fermi01 OpenMPI]$ mpirun -np 3 ./bcast
30 [0] Valeur diffuse : 4
31 [1] Valeur diffuse : 4
32 [2] Valeur diffuse : 4

```

2.2.11 MPI_Reduce

Syntaxe C

```

#include <mpi.h>
int MPI_Reduce(void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype,
MPI_Op op, int root, MPI_Comm comm)

```

Description

MPI_REDUCE effectue une opération de réduction (somme, max, ET logique ...) sur un ensemble de valeurs fournies par les processus du communicateur. La routine réduit, en utilisant l'opération `op`, les éléments contenus dans les buffers `sendbuf` et place le résultat dans le buffer `recvbuf` du processus de rang `root`. La taille des données se trouvant dans le buffer d'entrée et de sortie est définie par les arguments `count` et `datatype`. Tous les processus faisant appel MPI_REDUCE pour une même opération de réduction doivent donner des valeurs communes pour les arguments `count`, `datatype`, `op`, `root` et `comm`.

Cette opération est bloquante.

Paramètres d'entrée

`sendbuf` : adresse d'un des éléments qui subiront l'opération de réduction.

`count` : nombre d'éléments dans le buffer `sendbuf` d'un processus (différent du nombre total d'éléments qui seront réduits).

`datatype` : type des éléments réduits. Les types existants sont : MPI_CHAR, MPI_SHORT, MPI_INT, MPI_LONG, MPI_FLOAT, MPI_DOUBLE, MPI_LONG_LONG, MPI_LONG_DOUBLE⁹.

⁹Il est possible de créer des types de données dérivés, des types de données construits à partir des types fournis par MPI, à l'aide des routines MPI_Type_contiguous, MPI_Type_vector, MPI_Type_indexed, MPI_Type_struct,

op : opération de réduction à appliquer. Les opérations existantes sont : MPI_MAX (maximum), MPI_MIN (minimum), MPI_SUM (somme), MPI_PROD (produit), MPI_LAND (ET logique), MPI_BAND (ET bit--bit), MPI_LOR (OU logique), MPI_BOR (OU bit--bit), MPI_LXOR (OU exclusif logique), MPI_BXOR (OU exclusif bit--bit), MPI_MAXLOC et MPI_MINLOC (ces deux dernières opérations permettent de récupérer le maximum et le rang du processus qui le détient).

root : rang du processus qui recevra le résultat.

comm : commutateur.

Paramètres de sortie

recvbuf : adresse de la zone mémoire qui contiendra le résultat (doit être différente de sendbuf).

Valeur de retour

MPI_SUCCESS : la routine s'est exécutée sans erreur.

MPI_ERR_COMM : le commutateur est incorrect. Une erreur fréquente est d'utiliser un commutateur NULL.

MPI_ERR_COUNT : l'argument count est incorrect. L'argument ne peut pas être négatif.

MPI_ERR_TYPE : l'argument datatype est incorrect, le type n'est pas connu de MPI.

MPI_ERR_BUFFER : le pointeur buffer est incorrect. La zone de mémoire pointée doit être allouée. Cette erreur peut également apparaître s'il y a un recouvrement (aliasing)¹⁰ entre deux buffers.

MPI_ERR_ROOT : le rang du processus source est incorrect.

Illustration

Listing 2.11: Illustration de la fonction MPI_Reduce

```

1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, rank = 0;
7     MPI_Init(&argc, &argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD, &size);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11    MPI_Request sendState;
12
13    srand(rank);
14    int toreduce = rand()%10, reduced=0;
15    printf("[%d] Valeur : %d\n", rank, toreduce);
16    MPI_Reduce(&toreduce, &reduced, 1, MPI_INT, MPI_MAX, 0, MPI_COMM_WORLD);
17    if(rank==0)
18        printf("[%d] Valeur max : %d\n", rank, reduced);
19    MPI_Finalize();
20 }
21
22 [fremals@fermi01 OpenMPI]$ mpirun -np 5 reduce
23 [0] Valeur : 3
24 [4] Valeur : 1
25 [1] Valeur : 3
26 [3] Valeur : 6
27 [2] Valeur : 0

```

MPI_Type_commit, MPI_Type_free ...

¹⁰Le recouvrement est le recouvrement partiel ou total de deux zones mémoire. Ex. : une première zone mémoire pointe vers un tableau de 10 éléments tandis qu'une seconde zone mémoire pointe vers le 3^e élément du tableau.

```
28 [0] Valeur max : 6
```

2.2.12 MPI_Scatter

Syntaxe C

```
#include <mpi.h>
int MPI_Scatter(void *sendbuf, int sendcount, MPI_Datatype sendtype,
void *recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
```

Description

Le processus **root** distribue les données contenues dans **sendbuf** entre tous les processus du groupe dans l'ordre de leur rang. Chaque processus, le **root** y compris, aura ainsi **recvcount** éléments de type **recvtype** (dont la taille doit être identique à celle de **sendcount** éléments de type **sendtype**) dans leur buffer **recvbuf**. Les arguments **root** et **comm** doivent avoir la même valeur pour tous les processus appelant.

Cette opération est bloquante.

Paramètres d'entrée

sendbuf : adresse des éléments à distribuer (cet argument est significatif pour le processus **root** uniquement).

sendcount : nombre d'éléments à distribuer (cet argument est significatif pour le processus **root** uniquement).

sendtype : type des éléments à distribuer. Les types existants sont : **MPI_CHAR**, **MPI_SHORT**, **MPI_INT**, **MPI_LONG**, **MPI_FLOAT**, **MPI_DOUBLE**, **MPI_LONG**, **MPI_LONG_DOUBLE**¹¹ (cet argument est significatif pour le processus **root** uniquement).

recvcount : nombre d'éléments reçus.

recvtype : type des éléments distribués.

root : rang du processus qui distribue les données.

comm : communicateur.

Paramètres de sortie

recvbuf : adresse de la zone mémoire qui contiendra les éléments distribués.

Valeur de retour

MPI_SUCCESS : la routine s'est exécutée sans erreur.

MPI_ERR_COMM : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur **NULL**.

MPI_ERR_COUNT : l'argument **count** est incorrect. L'argument ne peut pas être négatif.

MPI_ERR_TYPE : l'argument **datatype** est incorrect, le type n'est pas connu de MPI.

MPI_ERR_BUFFER : le pointeur **buffer** est incorrect. La zone de mémoire pointée doit être allouée.

¹¹Il est possible de créer des types de données dérivés, des types de données construits à partir des types fournis par MPI, à l'aide des routines **MPI_Type_contiguous**, **MPI_Type_vector**, **MPI_Type_indexed**, **MPI_Type_struct**, **MPI_Type_commit**, **MPI_Type_free** ...

Illustration

Listing 2.12: Illustration de la fonction MPI_Scatter

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <mpi.h>
4
5 int main(int argc, char * argv[])
6 {
7     int size = 0, rank = 0, i;
8
9     MPI_Init(&argc,&argv);
10
11     MPI_Comm_size(MPI_COMM_WORLD,&size);
12     MPI_Comm_rank(MPI_COMM_WORLD,&rank);
13
14     int * data;
15     int myData;
16
17     if(rank == 0){
18         data = (int *) malloc(size*sizeof(int));
19         for(i=0; i<size; ++i){
20             data[i] = i;
21             printf("[%d] Data[%d] : %d\n", rank, i, data[i]);
22         }
23     }
24
25     MPI_Scatter(data, 1, MPI_INT, &myData, 1, MPI_INT, 0, MPI_COMM_WORLD);
26
27     ++ myData;
28     printf("[%d] Data : %d\n", rank, myData);
29
30
31     MPI_Gather(&myData, 1, MPI_INT, data, 1, MPI_INT, 0, MPI_COMM_WORLD);
32
33     if(rank==0){
34         for(i=0; i<size; ++i){
35             printf("[%d] Data[%d] : %d\n", rank, i, data[i]);
36         }
37     }
38
39     MPI_Finalize();
40 }
41
42 [fremals@fermi01 OpenMPI]$ mpirun -np 3 scatter
43 [0] Data[0] : 0
44 [0] Data[1] : 1
45 [0] Data[2] : 2
46 [0] Data : 1
47 [0] Data[0] : 1
48 [0] Data[1] : 2
49 [0] Data[2] : 3
50 [1] Data : 2
51 [2] Data : 3
```

2.2.13 MPI_Gather

Syntaxe C

```
#include <mpi.h>
int MPI_Gather(void *sendbuf, int sendcount, MPI_Datatype sendtype,
void *recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)
```

Description

Cette opération est le dual de `MPI_Scatter`. Chaque processus (le processus de rang `root` inclus) envoie le contenu du buffer `sendbuf` au processus de rang `root`. Les données sont concaténées dans le buffer de réception `recvbuf` dans l'ordre des rangs des processus. Les arguments `sendcount` et `sendtype` du processus `i` doivent indiquer une taille identique celle indiquée par les arguments `recvcount` et `recvtype` du processus `root` (`recvcount` indique donc le nombre d'éléments reçus de la part d'un processus, et non le nombre total d'éléments reçus). Les arguments `root` et `comm` doivent avoir la même valeur pour tous les processus appelant.

Cette opération est bloquante.

Paramètres d'entrée

`sendbuf` : adresse des éléments à regrouper.

`sendcount` : nombre d'éléments à transférer au `root`.

`sendtype` : type des éléments à regrouper. Les types existants sont : `MPI_CHAR`, `MPI_SHORT`, `MPI_INT`, `MPI_LONG`, `MPI_FLOAT`, `MPI_DOUBLE`, `MPI_LONG`, `MPI_LONG_DOUBLE`¹².

`recvcount` : nombre d'éléments envoyés par chaque processus et reçus par le `root` (cet argument est significatif pour le processus `root` uniquement).

`recvtype` : type des éléments regroupés (cet argument est significatif pour le processus `root` uniquement).

`root` : rang du processus qui reçoit les données.

`comm` : communicateur.

Paramètres de sortie

`recvbuf` : adresse de la zone mémoire qui contiendra les éléments concaténés (seul le processus `root` doit fournir une adresse valide).

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

`MPI_ERR_COMM` : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur `NULL`.

`MPI_ERR_COUNT` : l'argument `count` est incorrect. L'argument ne peut pas être négatif.

`MPI_ERR_TYPE` : l'argument `datatype` est incorrect, le type n'est pas connu de MPI.

`MPI_ERR_BUFFER` : le pointeur `buffer` est incorrect. La zone de mémoire pointée doit être allouée.

2.2.14 MPI_Wait

Syntaxe C

```
#include <mpi.h>
int MPI_Wait(MPI_Request *request, MPI_Status *status)
```

¹²Il est possible de créer des types de données dérivés, des types de données construits à partir des types fournis par MPI, à l'aide des routines `MPI_Type_contiguous`, `MPI_Type_vector`, `MPI_Type_indexed`, `MPI_Type_struct`, `MPI_Type_commit`, `MPI_Type_free` ...

Description

Un appel `MPI_Wait` rend la main quand l'opération identifiée par la variable `request` est achevée. L'objet représentant la requête de communication est dsalloué par cette routine et la variable `request` est mise à la valeur `MPI_REQUEST_NULL`.

La variable `status` contient des informations sur l'opération achevée. Le champ `MPI_ERROR` de cette variable contient des informations lorsque `MPI_Testall`, `MPI_Testsome`, `MPI_Waitall` ou `MPI_Waitsome` a renvoyé l'erreur `MPI_ERR_IN_STATUS`. Dans le cas des routines d'envoi, `status` est utile uniquement pour la fonction `MPI_Test_cancelled` ou si l'erreur `MPI_ERR_STATUS` a été retournée. La valeur du champ `MPI_ERROR` vaut alors `MPI_SUCCESS` pour chaque opération d'envoi ou de réception réussie, ou `MPI_ERR_PENDING` pour les opérations qui ont échoué ou qui sont inachevées. S'il n'est pas nécessaire d'examiner la variable `status`, celle-ci peut être remplacée par la constante `MPI_STATUS_IGNORE`.

Cette opération est bloquante.

Paramètres d'entrée

`request` : identifie l'opération à attendre.

Paramètres de sortie

`status` : structure contenant des informations sur l'opération ; le paramètre peut être remplacé par `MPI_STATUS_IGNORE`.

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

`MPI_ERR_REQUEST` : l'argument `request` est incorrect.

`MPI_ERR_ARG` : un argument est incorrect.

Illustration

Listing 2.13: Illustration de la fonction `MPI_Wait`

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int computation(){
5     return rand()%10;
6 }
7
8 int main(int argc, char * argv[])
9 {
10     int size = 0, rank = 0;
11     MPI_Init(&argc,&argv);
12
13     MPI_Comm_size(MPI_COMM_WORLD,&size);
14     MPI_Comm_rank(MPI_COMM_WORLD,&rank);
15     MPI_Request sendState;
16     if(rank == 0){
17         sleep(1);
18         int tosend = computation(), temp;
19         MPI_Isend(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, &sendState);
20         while(tosend!=9){
21             temp = computation(); //le buffer d'envoi ne peut être modifié
22             MPI_Wait(&sendState, MPI_STATUS_IGNORE);
```

```

23         tosend = temp;
24         MPI_Isend(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, &sendState);
25     }
26     MPI_Wait(&sendState, MPI_STATUS_IGNORE);
27 }
28 if(rank==1){
29     MPI_Status status;
30     int torecv=0;
31     while(torecv!=9){
32         MPI_Recv(&torecv, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
33         printf("Valeur calcule : %d\n", torecv);
34     }
35 }
36 MPI_Finalize();
37 }
38
39 [fremals@fermi01 OpenMPI]$ mpirun -np 2 wait
40 Valeur calcule : 3
41 Valeur calcule : 6
42 Valeur calcule : 7
43 Valeur calcule : 9

```

2.2.15 MPI_Test

Syntaxe C

```

#include <mpi.h>
int MPI_Test(MPI_Request *request, int *flag, MPI_Status *status)

```

Description

La routine renseigne sur l'état d'une opération identifiée par `request`. `flag` contient 1 si l'opération s'est terminée et 0 sinon.

La variable `status` contient des informations sur l'opération achevée. Le champ `MPI_ERROR` de cette variable contient des informations lorsque `MPI_Testall`, `MPI_Testsome`, `MPI_Waitall` ou `MPI_Waitsome` a renvoyé l'erreur `MPI_ERR_IN_STATUS`. Dans le cas des routines d'envoi, `status` est utile uniquement pour la fonction `MPI_Test_cancelled` ou si l'erreur `MPI_ERR_STATUS` a été retournée. La valeur du champ `MPI_ERROR` vaut alors `MPI_SUCCESS` pour chaque opération d'envoi ou de réception réussie, ou `MPI_ERR_PENDING` pour les opérations qui ont échoué ou qui sont incomplètes. S'il n'est pas nécessaire d'examiner la variable `status`, celle-ci peut être remplacée par la constante `MPI_STATUS_IGNORE`.

Cette opération est non-bloquante.

Paramètres d'entrée

`request` : identifie l'opération sondée.

Paramètres de sortie

`flag` : indique si l'opération s'est achevée (=1) ou non (=0).

`status` : structure contenant des informations sur l'opération ; le paramètre peut être remplacé par `MPI_STATUS_IGNORE`.

Valeur de retour

MPI_SUCCESS : la routine s'est exécutée sans erreur.

MPI_ERR_REQUEST : l'argument `request` est incorrect.

MPI_ERR_ARG : un argument est incorrect.

Illustration

Listing 2.14: Illustration de la fonction `MPI_Test`

```

1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, rank = 0;
7     MPI_Init(&argc, &argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD, &size);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11    if(rank == 0){
12        sleep(1);
13        int tosend = 5;
14        MPI_Send(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
15    }
16    if(rank==1){
17        MPI_Status status;
18        int torecv;
19        MPI_Request recvState;
20        MPI_Irecv(&torecv, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &recvState);
21        int recved = 0;
22        while(!recved){
23            MPI_Test(&recvState, &recved, MPI_STATUS_IGNORE);
24            printf("Etat de la r ception : %d\n", recved);
25            if(!recved)
26                sleep(0);
27        }
28    }
29    MPI_Finalize();
30 }
31
32 [fremals@fermi01 OpenMPI]$ mpirun -np 2 ./irecv
33 Etat de la r ception : 0
34 Etat de la r ception : 0
35 ...
36 Etat de la r ception : 0
37 Etat de la r ception : 0
38 Etat de la r ception : 1

```

2.2.16 MPI_Probe

Syntaxe C

```

#include <mpi.h>
int MPI_Probe(int source, int tag, MPI_Comm comm, MPI_Status *status)

```

Description

L'appel `MPI_Probe` est bloquant jusqu' l'envoi par le processus de rang `source` d'un message portant l'étiquette `tag` dans le communicateur `comm`. Le message n'est pas reçu, mais cela permet l'utilisateur de décider comment le recevoir en se basant sur les informations présentes dans `status`. Il peut, par exemple, allouer de la mémoire pour le buffer de réception en fonction de la longueur du message (`MPI_Get_Count` permet de connaître le nombre d'éléments envoyés).

Si le champ `status` ne doit pas être examiné, il peut être remplacé par la constante `MPI_STATUS_IGNORE`.

Paramètres d'entrée

`source` : rang du processus source ou `MPI_ANY_SOURCE`.

`tag` : valeur de l'étiquette du message attendu ou `MPI_ANY_TAG`. La plus grande valeur possible pour ce paramètre est contenue dans la constante `MPI_TAG_UB`.

`comm` : communicateur.

Paramètres de sortie

`status` : structure contenant des informations sur le message sondé.

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

`MPI_ERR_COMM` : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur `NULL`.

`MPI_ERR_TAG` : l'argument `tag` est incorrect. Les tags doivent être non-négatifs ; la plus grande valeur permise pour ce paramètre est disponible dans l'attribut `MPI_TAG_UB`. Dans le cas de la réception, l'étiquette peut valoir `MPI_ANY_TAG`.

`MPI_ERR_RANK` : rang du destinataire incorrect. Le rang doit être compris entre 0 et la taille du groupe du communicateur décrement d'une unité. Dans le cas de la réception, le rang peut valoir `MPI_ANY_SOURCE`.

Illustration

Listing 2.15: Illustration de la fonction `MPI_Probe`

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <mpi.h>
4
5 int main(int argc, char * argv[])
6 {
7     srand(5);
8     int size = 0, rank = 0;
9     MPI_Init(&argc, &argv);
10
11     MPI_Comm_size(MPI_COMM_WORLD, &size);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
13     if(rank == 0){
14         MPI_Status status;
15         int compute1=0;
16         double compute2=0;
17         while(compute1!=9 || compute2<8){
18             MPI_Probe(MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
19             switch(status.MPI_TAG){

```

```

20         case 1 :
21             MPI_Recv(&compute1, 1, MPI_INT, status.MPI_SOURCE,
22                     status.MPI_TAG, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
23             printf("Premier calcul : %d\n", compute1);
24             break;
25         case 2 :
26             MPI_Recv(&compute2, 1, MPI_DOUBLE, status.MPI_SOURCE,
27                     status.MPI_TAG, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
28             printf("Deuxime calcul : %f\n", compute2);
29             break;
30     }
31 }
32 }
33 else if(rank==1){
34     int compute1 = 0;
35     while(compute1!=9){
36         compute1 = computation1();
37         printf("Compute1 : %d\n", compute1);
38         MPI_Send(&compute1, 1, MPI_INT, 0, 1, MPI_COMM_WORLD);
39     }
40 }
41 else if(rank==2){
42     double compute2 = 0;
43     while(compute2<8){
44         compute2 = computation2();
45         printf("Compute2 : %f\n", compute2);
46         MPI_Send(&compute2, 1, MPI_DOUBLE, 0, 2, MPI_COMM_WORLD);
47     }
48 }
49 MPI_Finalize();
50 }
51
52 [fremals@fermi01 Listings]$ mpirun -np 3 ./probe
53 Deuxime calcul : 2.747456
54 Deuxime calcul : 0.464678
55 Deuxime calcul : 9.927552
56 Compute2 : 2.747456
57 Compute2 : 0.464678
58 Compute2 : 9.927552
59 Compute1 : 5
60 Compute1 : 5
61 ...
62 Compute1 : 0
63 Compute1 : 9
64 Premier calcul : 5
65 Premier calcul : 5
66 ...
67 Premier calcul : 0
68 Premier calcul : 9

```

2.2.17 MPI_Iprobe

Syntaxe C

```

#include <mpi.h>
int MPI_Iprobe(int source, int tag, MPI_Comm comm, int *flag,
MPI_Status *status)

```

Description

L'appel `MPI_Iprobe` permet de vérifier sans blocage la présence d'un message entrant qui porte l'étiquette `tag` et qui est mis par le processus de rang `source` dans le communicateur `comm`. Le paramètre `flag` est mis à 1 si un message peut être reçu ; le paramètre `status` contient alors les mêmes informations qu'aurait fourni un appel `MPI_Recv()`. Le message n'est pas reçu, mais cela permet à l'utilisateur de décider comment le recevoir en se basant sur les informations présentes dans `status` (il peut, par exemple, allouer de la mémoire pour le buffer de réception en fonction de la longueur du message).

Si le champ `status` ne doit pas être examiné, il peut être remplacé par la constante `MPI_STATUS_IGNORE`.

Paramètres d'entrée

`source` : rang du processus source ou `MPI_ANY_SOURCE`.

`tag` : valeur de l'étiquette du message attendu ou `MPI_ANY_TAG`. La plus grande valeur possible pour ce paramètre est contenue dans la constante `MPI_TAG_UB`.

`comm` : communicateur.

Paramètres de sortie

`status` : structure contenant des informations sur le message sondé.

`flag` : booléen indiquant la présence ou non d'un message entrant correspondant aux critères indiqués.

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

`MPI_ERR_COMM` : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur `NULL`.

`MPI_ERR_TAG` : l'argument `tag` est incorrect. Les tags doivent être non-négatifs.

`MPI_ERR_RANK` : rang du destinataire incorrect. Le rang doit être compris entre 0 et la taille du groupe du communicateur décimé d'une unité. Dans le cas de la réception, le rang peut valoir `MPI_ANY_SOURCE`.

Illustration

Listing 2.16: Illustration de la fonction `MPI_Iprobe`

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main( int argc, char * argv[] )
5 {
6     int rank;
7     int sendMsg = 123;
8     int recvMsg = 0;
9     int flag = 0;
10    int count;
11    MPI_Status status;
12    MPI_Request request;
13    int errs = 0;
14
15    MPI_Init(&argc, &argv);
16    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
```

```
17
18     if(rank == 0)
19     {
20         MPI_Isend(&sendMsg, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, &request);
21         while(!flag)
22         {
23             MPI_Iprobe(0, 0, MPI_COMM_WORLD, &flag, &status);
24         }
25         MPI_Get_count(&status, MPI_INT, &count);
26         if(count != 1)
27         {
28             errs++;
29         }
30         MPI_Recv(&recvMsg, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, &status);
31         if (recvMsg != 123)
32         {
33             errs++;
34         }
35         MPI_Wait(&request, &status);
36     }
37
38     MPI_Finalize();
39     return errs;
40 }
```

2.2.18 MPI_Get_count

Syntaxe C

```
#include <mpi.h>
int MPI_Get_count( MPI_Status *status, MPI_Datatype datatype, int *count )
```

Description

Cette routine permet de connaître le nombre d'éléments en cours d'envoi pour un transfert identifié par la variable `status`.

Paramètres d'entrée

`status` : le statut retourné par une opération de réception ou de sondage.
`datatype` : le type des éléments reçus.

Paramètres de sortie

`count` : le nombre d'éléments qui seront reçus. Si la taille de `datatype` vaut 0, `count` est mis à 0 ; si la taille des données dans `status` n'est pas un multiple de la taille de `datatype`, `count` vaut `MPI_UNDEFINED`.

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

Illustration

Listing 2.17: Illustration de la fonction `MPI_Get_count`

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <mpi.h>
4
5 int computation1(){
6     return rand()%10;
7 }
8
9 double computation2(){
10     return (double) rand()/RAND_MAX*10;
11 }
12
13 int main(int argc, char * argv[])
14 {
15     srand(5);
16     int size = 0, rank = 0;
17     MPI_Init(&argc,&argv);
18
19     MPI_Comm_size(MPI_COMM_WORLD,&size);
20     MPI_Comm_rank(MPI_COMM_WORLD,&rank);
21     if(rank == 0){
22         MPI_Status status;
23         int count;
24
25         MPI_Probe(MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
26         MPI_Get_count(&status, MPI_INT, &count);
27         printf("[%d] Taille du message : %d\n", rank, count);
28
29         int * numbers = (int *) malloc(count*sizeof(int));
30         MPI_Recv(numbers, count, MPI_INT, status.MPI_SOURCE, status.MPI_TAG,
31                 MPI_COMM_WORLD, MPI_STATUS_IGNORE);
32         free(numbers);
33     }
34     else if(rank==1){
35         int len = rand()%10, i;
36         int * numbers = (int *) malloc(len*sizeof(int));
37         for(i=0; i<len;++i)
38             numbers[i] = i;
39         printf("[%d] Taille du message : %d\n", rank, len);
40         MPI_Send(numbers, len, MPI_INT, 0, 1, MPI_COMM_WORLD);
41         free(numbers);
42     }
43     MPI_Finalize();
44 }
45
46 [fremals@fermi01 OpenMPI]$ mpirun -np 2 get_count
47 [1] Taille du message : 5
48 [0] Taille du message : 5

```

2.2.19 MPI_Barrier

Syntaxe C

```

#include <mpi.h>
int MPI_Barrier(MPI_Comm comm)

```

Description

L'appel `MPI_Barrier` bloque l'exécution du processus appelant jusqu'à ce que tous les processus du communicateur aient atteint cette routine.

Paramètres d'entrée

`comm` : communicateur.

Valeur de retour

`MPI_SUCCESS` : la routine s'est exécutée sans erreur.

`MPI_ERR_COMM` : le communicateur est incorrect. Une erreur fréquente est d'utiliser un communicateur `NULL`.

2.2.20 MPI_Wtime

Syntaxe C

```
#include <mpi.h>
double MPI_Wtime();
```

Description

Cette fonction renvoie le nombre flottant de secondes écoulées depuis un certain moment dans le passé. Pour mesurer un laps de temps, il est donc nécessaire de faire deux mesures et de soustraire la première de la seconde.

Valeur de retour

Retourne le nombre flottant de secondes écoulées depuis un certain point dans le passé.

Illustration

Listing 2.18: Illustration de la fonction `MPI_Wtime`

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc, char * argv[])
5 {
6     int size = 0, rank = 0;
7     MPI_Init(&argc,&argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD,&size);
10    MPI_Comm_rank(MPI_COMM_WORLD,&rank);
11    if(rank == 0){
12        int tosend = 5;
13        MPI_Ssend(&tosend, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);
14        printf("[%d] Message post .\n", rank);
15    }
16    if(rank==1){
17        double start, end;
18        start = MPI_Wtime();
19        sleep(1);
20        MPI_Status status;
21        int torecv;
22        printf("[%d] Lecture de la boîte au lettre.\n", rank);
23        MPI_Recv(&torecv, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
24        end = MPI_Wtime();
25        printf("[%d] Message reçu : %d (%f s écoulées).\n", rank, torecv,
26                end-start);
27    }
28    MPI_Finalize();
```



```
29 }  
30  
31 [fremals@fermi01 OpenMPI]$ mpirun -np 2 wtime  
32 [0] Message post .  
33 [1] Lecture de la boite au lettre.  
34 [1] Message reçu : 5 (1.000319 secondes ).
```



Les processeurs graphiques et CUDA

3.1 Introduction

Les ordinateurs actuels présentent une architecture complexe pouvant associer des processeurs centraux multi-cœurs (CPU - Central Processing Unit multicore), ainsi que des processeurs graphiques (GPU - Graphics Processing Unit) de plus en plus aptes exécuter des calculs généralistes. Ce paragraphe est présent en trois parties:

1. Architecture des processeurs centraux et graphiques.
2. Mémoires des processeurs graphiques.
3. Langages de programmation sur GPU.
4. Outils d'exploitation des architectures hétérogènes (multi-CPU/multi-GPU).

3.2 Architecture des processeurs centraux et graphiques



Figure 3.1: Présentation des architectures des CPUs et des GPUs

Les processeurs centraux et les processeurs graphiques¹ sont représentés la figure 3.1. Le processeur central est composé d'une Unité Arithmétique et Logique (ALU, Arithmetic and Logical Unit) qui effectue principalement les calculs mathématiques (addition, soustraction, multiplication et division) et logiques (OR, AND, NOT, XOR, tests binaires ...) ; d'une unité de contrôle, ou séquenceur, qui synchronise les différents éléments du processeur (il charge les instructions du

¹Les processeurs graphiques dont nous parlons dans ce document correspondent aux modèles développés par NVIDIA. Les processeurs graphiques AMD/ATI peuvent contenir des éléments qui diffèrent.

programme en mémoire, les decode et les fait exécuter par les ALUs de manière séquentielle) ; et de plusieurs niveaux de mémoire cache qui permet l'accélération des traitements (mémoire rapide mais coûteuse, proche du CPU). Enfin, la DRAM est la mémoire vive où les données sont placées lorsqu'elles sont traitées. Cette mémoire ne se trouve pas sur le processeur mais au niveau de la carte mère, sous forme de barrette mémoire.

Pour les CPU, nous appelons "cœur processeur" la combinaison d'une ALU et d'un contrôleur. Plus il y en a, et plus il est possible d'effectuer des calculs simultanément².

Les processeurs graphiques sont découps en multi-processeurs (SM, Streaming Multiprocessor), chaque multi-processeur est composé de 8, 16 ou 32 cœurs (SM, Streaming Processor). Les cœurs des GPUs sont moins performants que ceux des CPUs, mais leur grand nombre (128 dans l'exemple) permet le traitement simultané d'un grand ensemble de données et de rendre ainsi le GPU plus performant que le CPU. Chaque SM est ici un processeur composé de 8, 16 ou 32 ALUs, d'une unité de contrôle et d'une zone de mémoire cache. Les ALUs des GPUs sont généralement moins performantes que celles des CPUs (vitesse d'horloge³). Cette structure matérielle particulière a été développée pour maximiser les traitements graphiques typiques des cartes vidéos, où les mêmes traitements sont simultanément appliqués sur des données différentes (mode de fonctionnement SIMD, Single Instruction Multiple Data). Ainsi, toutes les ALUs d'un même SM effectuent la même instruction au même moment.

Les GPU nouvelle génération (Fermi) apportent des nouvelles améliorations, notamment les SM indépendants. Des programmes différents peuvent s'exécuter concurremment sur les SMs.

3.3 Structures des mémoires dédiées aux processeurs graphiques

La structure mémoire dédiée au processeur graphique est appelée mémoire device. Elle est physiquement située sur la carte mémoire proche du processeur graphique. Un ordinateur peut posséder plusieurs cartes d'extension, ou device, dont les cartes graphiques font parties. Par opposition, la mémoire dédiée aux CPUs est appelée mémoire hôte car c'est la mémoire située sur la carte mère de l'ordinateur; celle qui héberge les cartes d'extension.

La structure mémoire dédiée aux processeurs graphiques est complexe mais importante à maîtriser. En effet, c'est l'utilisateur qui a la charge du placement des données dans les différents types de la mémoire device. Il y a cinq types de mémoire (cf. fig. 3.2) :

1. Mémoire globale : correspond à la mémoire RAM du processeur graphique. C'est l'un des types de mémoire accessible par le CPU. La capacité de cette mémoire varie généralement de 512 Mo à 4 Go en fonction du type de la carte graphique. Cette mémoire de masse se trouvant sur la carte graphique a une latence d'accès aux données importante : de 400 à 600 cycles d'horloge.
2. Mémoire partagée : cette mémoire se trouve dans le SM. Une donnée placée dans la mémoire partagée d'un SM n'est accessible que par n'importe lequel de ses cœurs. Cette mémoire est limitée à 16 ko par SM et possède une latence de 4 cycles d'horloge. Si une donnée est accédée en lecture par tous les threads tournant sur le SM, un système de broadcast matériel permet une propagation de la donnée sans qu'il y ait de conflit d'accès. Il y aura par contre un conflit d'accès si certains threads, mais pas tous, accèdent à une même donnée.
3. Mémoire constante : cette mémoire est localisée dans le SM. Ce type de mémoire est accessible par le CPU en écriture mais uniquement en lecture par le GPU. Elle est physiquement optimisée pour le cas où tous les threads s'exécutent accèdent à la même donnée (par exemple,

²Par défaut, le parallélisme est obtenu dans le cadre de ce cours en utilisant les threads.

³Environ 1,5 GHz pour les GPUs et environ 3 GHz pour les CPUs

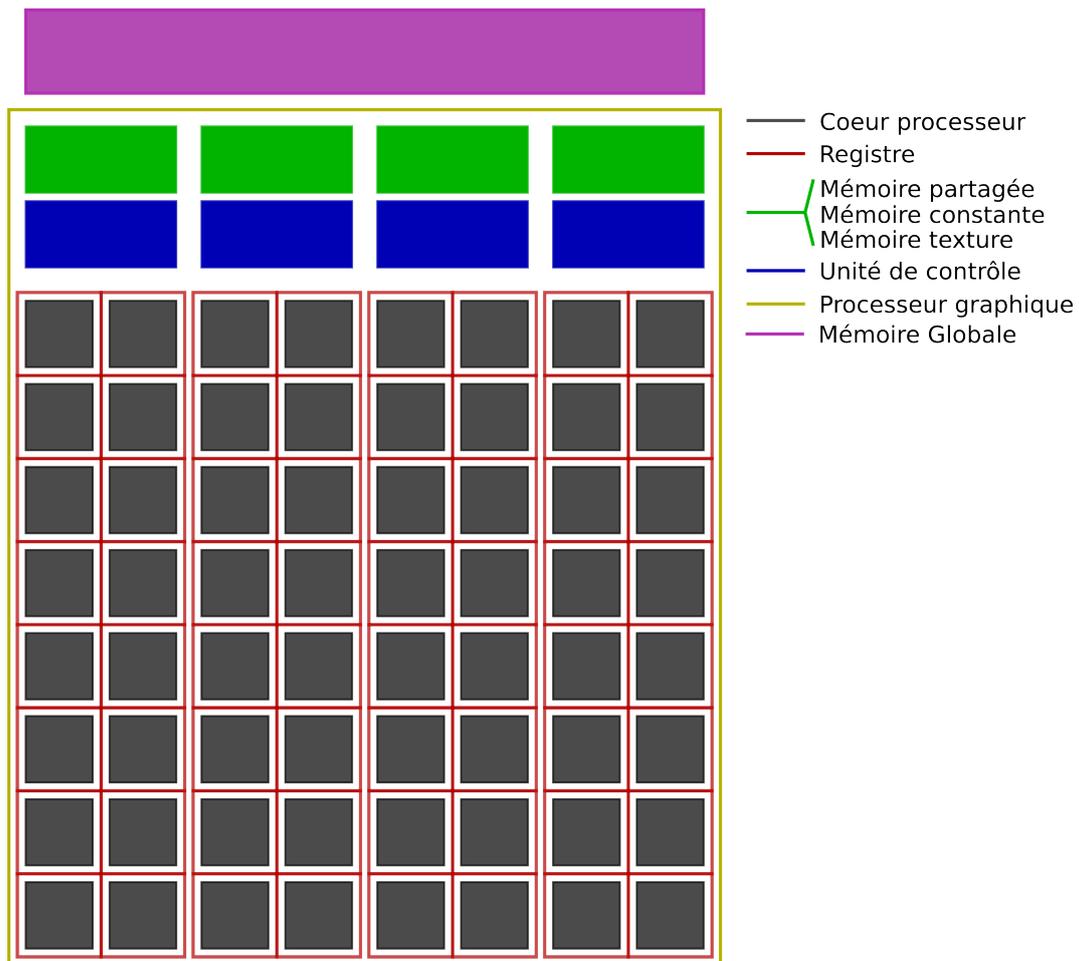


Figure 3.2: Présentation de la structure des mémoires d'un GPU de 64 cœurs

c'est le cas de la constante π), sa latence est alors de 1 cycle. Lorsque les threads accèdent simultanément des adresses différentes de la mémoire constante, les accès sont sérialisés⁴. La capacité de cette mémoire est limitée à 64 ko par GPU et 8 ko par SM.

4. Mémoire texture: tout comme la mémoire constante, cette mémoire est en lecture seule pour le GPU et accessible en écriture par le CPU. Le coût de la lecture est aussi très faible (4 cycles d'horloge). Les textures permettent vraiment de simplifier le traitement d'images : elles permettent la mise en œuvre de filtrages bilinéaires et trilineaires très facilement et l'accès aléatoire aux pixels.
5. Registre : cette mémoire correspond aux registres de travail des ALUs, ce qui en fait la mémoire qui possède la latence la plus basse (1 cycle d'horloge). Il y a 8192 registres de quatre octets par SM.
6. Mémoire locale : cette mémoire, de même type que la mémoire globale (latence de 400 à 600 cycles d'horloge) est utilisée comme back up, ou swap, de registres lorsque le nombre de registres est insuffisant pour un programme. Elle a le même rôle que le swap disk pour les CPUs. Sa taille est déterminée à la compilation. En général, son emploi est évité pour éviter les mauvaises performances.

La Table 3.1 décrit les propriétés des différentes mémoires du GPU. Les tailles données sont des

⁴Effectués les uns après les autres.

ordres de grandeur, car elles peuvent varier selon le modèle de carte graphique, le nombre de registres disponibles par thread dépend de la programmation de l'algorithme.

Type de mémoire	Utilité	Taille	Latence (cycles)
Registres	Propre à chaque thread	8192×4 octets	1
Locale	Complète les registres	Indéterminée	400 - 600
Partage	Communication entre threads	16 ko	4
Constante	Lecture seule	8 ko	1 ou plus
Textures	Utile en traitement d'images	1 ko par unité de calcul	4
Globale	Mémoire principale	jusqu'à 4 Go	400 - 600

Table 3.1: Propriétés des mémoires du GPU

3.4 Les langages de programmation GPU

Comme montré dans l'introduction, les GPUs ont une architecture massivement parallèle. La possibilité de pouvoir les utiliser pour faire autre chose que le rendu d'images a attiré l'attention de nombreux chercheurs. Cependant, jusqu'au début des années 2000, les GPU n'étaient pas facilement programmables. Ils ne pouvaient qu'exécuter une suite fixe de traitements appliqués aux objets géométriques pour obtenir le rendu 2D ou 3D, cette suite de traitements appelée "*pipeline graphique*" était figée. Seules quelques expériences marginales utilisant certaines fonctions natives du matériel furent menées et permirent d'accélérer quelques tâches simples avec ce qu'on appelle les Shaders.

Les Shaders sont des programmes qui permettent de paramétrer une partie du pipeline graphique. Apparus à la fin des années 90, ils ont commencé à être beaucoup plus utilisés à partir de 2001. Depuis, ils n'ont pas cessé d'évoluer et d'offrir de plus en plus de possibilités aux programmeurs grâce à des langages permettant la programmation de ces Shaders. OpenGL et Direct3D sont deux Shaders très populaires, de plus en plus utilisés dans la création des applications de génération d'images 3D.

Les dernières années ont été marquées par l'ouverture à la programmation des GPUs avec l'avènement de nouveaux langages de programmation. CUDA et OpenCL sont les plus utilisés actuellement ; ils permettent de porter des algorithmes génériques sur les processeurs graphiques (GPGPU - General-Purpose computing on Graphics Processing Units). Nous présentons ci-dessous une description de la bibliothèque graphique OpenGL (rendu d'images 2D/3D) et des environnements de programmation sur GPU (CUDA et OpenCL). Notons qu'une description plus approfondie de CUDA est présente puisqu'il représente le langage de programmation GPU le plus utilisé actuellement, en raison de son efficacité et de sa compatibilité avec les cartes les plus utilisées : NVIDIA.

3.4.1 OpenGL

OpenGL (Open Graphic Library) [8] est une bibliothèque graphique ouverte, reconnue comme un standard et développée depuis 1992 par Silicon Graphics [9], un des titans de l'informatique graphique. Il a eu ses lettres de noblesse dans l'imagerie scientifique et dans la production

cinématographique. L'aspect ouvert d'OpenGL lui a permis d'être port sur des plateformes et systèmes d'exploitation variés. Son concurrent direct est la bibliothèque Direct 3D développée par Microsoft et fortement reconnue dans le monde du jeu vidéo.

L'objectif majeur d'OpenGL est de construire et de visualiser des scènes 2D/3D (images de synthèse, photos, etc.) sur écran. OpenGL permet l'interactivité avec l'utilisateur en gérant les périphériques d'entrée (clavier, souris, etc.) pour modifier la scène (déplacement des objets, changement des points de vue de la caméra, etc.) et de visualiser l'effet immédiatement. OpenGL est préconisé dans les milieux scientifiques et académiques du fait de son ouverture, de sa souplesse d'utilisation et de sa disponibilité sur des plateformes variées.

3.4.2 CUDA : programmation des processeurs graphiques NVIDIA

L'API CUDA (Compute Unified Device Architecture) représente l'une des technologies les plus utilisées pour la programmation parallèle sur GPU. CUDA a été développée par NVIDIA en 2007, permettant de programmer les processeurs graphiques dans un langage proche du C standard⁵

Utiliser un processeur graphique pour exécuter un programme sous CUDA est un procédé requérant plusieurs étapes (cf. figure 3.3). La première étape est facultative et consiste à sélectionner le processeur graphique qui sera utilisé pour le traitement. Cette étape est utile uniquement si la machine possède plusieurs GPUs et que l'utilisateur veut soit les utiliser en même temps, soit utiliser un autre GPU que celui assigné par défaut. La seconde étape est de réserver l'espace nécessaire, au niveau de la mémoire globale du device, pour toute donnée d'entrée et/ou de sortie. Les données d'entrée sont ensuite transférées de la mémoire hôte vers la mémoire device allouée, les rendant ainsi disponibles pour leur traitement sur le processeur graphique. L'utilisateur doit ensuite spécifier le nombre de threads à lancer sur le GPU ainsi que leur agencement. Les traitements sont alors exécutés par les différents cœurs du processeur graphique. Enfin, les résultats sont rapatriés sur la mémoire hôte et l'espace mémoire device alloué peut être libre.

Chacune de ces étapes demande des connaissances propres aux processeurs graphiques et au langage de programmation CUDA. Celles-ci sont présentées en détail dans les sept points suivants:

1. **Sélection du processeur graphique:** Un processeur graphique est par défaut assigné pour effectuer les traitements requis par le processeur central, mais l'utilisateur peut choisir de changer de processeur s'il en possède plusieurs. Pour ce faire, il faut utiliser les fonctions `cudaGetDeviceCount` et `cudaSetDevice` :

```
cudaError_t cudaGetDeviceCount(int * count)
```

```
cudaError_t cudaSetDevice(int device)
```

La première fonction place le nombre de processeurs graphiques présents dans la machine au sein de la variable `count`. La seconde fonction permet de sélectionner le GPU et la variable `device` représente l'identifiant du processeur qui sera assigné pour les calculs. Celui-ci est compris entre 0 et `count-1`.

2. **Allocation de la mémoire:** Avant toute chose, il faut allouer la mémoire device nécessaire aux traitements sur le GPU (au minimum pour les données d'entrée et de sortie). Pour ce faire, CUDA implémente sa propre version de la fonction `malloc`⁶ :

```
cudaError_t cudaMalloc (void ** devPtr, size_t size)
```

⁵NVIDIA CUDA. <https://developer.nvidia.com/cuda-zone>

⁶La fonction `malloc` permet de réserver, sur la mémoire hôte, un espace mémoire d'une taille spécifique par le développeur

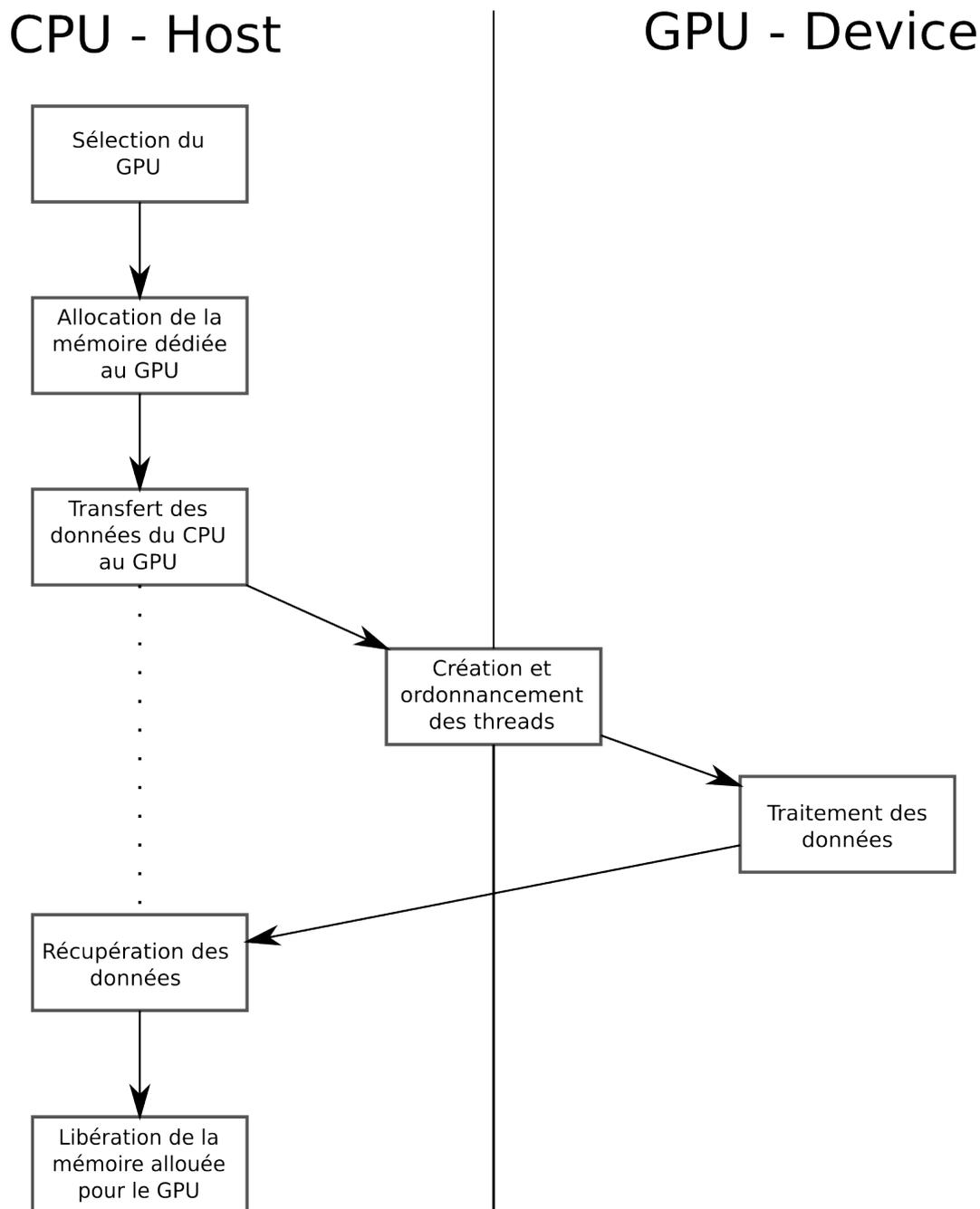


Figure 3.3: Etapes ncessaire un traitement sur processeur graphique

La variable `devPtr` est un pointeur qui recevra l'adresse du dbut de la zone mmoire device alloué et `size` dtermine la taille demande de l'espace mmoire.

- 3. Transfert des donnes CPU ⇒ GPU:** Pour effectuer le transfert des donnes entre la mmoire hte et la mmoire device, CUDA implmente sa propre version de la fonction `memcpy`⁷ :

```

cudaError_t cudaMemcpy(void * dst, const void * src, size_t count, enum
                      cudaMemcpyKind kind)
  
```

⁷La fonction `memcpy` permet de copier, sur le processeur graphique, les donnes d'un espace mmoire dans un autre



La variable `dst` est un pointeur vers l'espace mémoire qui accueillera les données, `src` est un pointeur vers l'espace mémoire qui contient les données, `count` indique la taille en octet des données transférer (l'espace mémoire de `src` `src+count` seront copiés vers l'espace mémoire `dst` `dst+count`), et `kind` indique le type de transfert à effectuer. `kind` connaît quatre valeurs :

- (a) `cudaMemcpyHostToHost` : la donnée est copiée depuis un espace mémoire hôte vers un espace mémoire hôte (cette fonction correspond à un appel `memcpy`)
 - (b) `cudaMemcpyHostToDevice` : la donnée est copiée depuis un espace mémoire hôte vers un espace mémoire device (c'est le type qui est utilisé dans le cadre de ce paragraphe)
 - (c) `cudaMemcpyDeviceToHost` : la donnée est copiée depuis un espace mémoire device vers un espace mémoire hôte
 - (d) `cudaMemcpyDeviceToDevice` : la donnée est copiée depuis un espace mémoire device vers un espace mémoire device
4. **L'appel du kernel : le partitionnement des threads:** Le GPU fonctionne selon le mode SIMD et est ainsi conçu pour que tous les cœurs exécutent simultanément les mêmes instructions. La fonction définissant la suite de ces instructions est appelée "kernel".

Le processeur graphique fonctionne uniquement avec des threads, il n'y a pas de notion de processus. Lors du lancement du kernel, le développeur indique le nombre de threads qui seront créés (valeur constante, figée lors de la compilation) et tous ces threads exécuteront ce même kernel. Les contrôleurs des multiprocesseurs gèrent les threads par paquets (warp⁸) de 32 :

- Pour les GPUs possédant des SM de 8 cœurs : un warp est actif à la fois, chaque cœur gère quatre threads
- Pour les GPUs possédant des SM de 32 cœurs : deux warps sont actifs à la fois, chaque cœur gère deux threads

Le contrôleur du SM gère les threads par paquet de 32. L'utilisateur peut également regrouper les threads au sein de blocs avec un maximum de 768 ou 1024 threads par blocs (selon le GPU, les autres permettent un maximum de 768). La répartition de la mémoire partagée est définie par la taille du bloc : les données qui sont stockées dans la mémoire partagée sont accessibles par tous les threads d'un même bloc. Lors de l'appel au kernel, l'utilisateur définit le nombre de blocs et le nombre de threads par blocs.

Lors de l'exécution, les blocs sont assignés aux SM de manière équivalente à l'assignation des processus sur les processeurs centraux par le système d'exploitation⁹. Les blocs sont décomposés en warp et ceux-ci sont exécutés par les SP. Lorsqu'un warp devient inactif (synchronisation, accès mémoire) le warp est interchangé avec un warp en attente d'exécution. Ce mécanisme permet de recouvrir les temps d'inactivité des warps par les temps d'activité d'autres warps. Son efficacité dépendra du nombre de threads, et donc de warps, assignés à un SM et de la durée des temps d'inactivité. Plus les temps d'inactivité sont longs, plus le nombre de threads doit être élevé pour obtenir une bonne couverture, et donc de meilleures performances.

Pour connaître le nombre de blocs assignés au SM, il suffit de résoudre ce système d'équations :

⁸Ce terme provient des métiers tisser, il désigne un ensemble de fils (thread = fil) ; un warp CUDA est ainsi un ensemble de 32 threads. [Source](#).

⁹Cf. algorithme de Round-Robin

- $p * b \leq 768, 1024$ ou 10536 : maximum 768, 1024 ou 1536 threads (selon le GPU) peuvent être attribués à un SM
- $r * p * b \leq R$: il ne faut pas dépasser la quantité de registres disponibles
- $s * b \leq S$: il ne faut pas dépasser la quantité de mémoire partagée disponible
- $b \leq 8$: limite matérielle
- $p \leq 512$: limite matérielle

Avec :

- R : l'ensemble des registres disponibles pour un SM
- r : le nombre de registres utilisés par un thread
- S : la quantité de mémoire partagée disponible dans un SM
- s : la quantité de mémoire partagée utilisée par un bloc
- b : le nombre de blocs assignés à un SM
- p : le nombre de threads par bloc

Dans ce système de contraintes, r et s sont fixés à la compilation du kernel et peuvent être fonction de p . b dépend également de p , c'est le nombre de threads par bloc qui permet d'ajuster la répartition des blocs sur les SM.

Listing 3.1: Appel d'un kernel : addition de vecteurs

```

1
2 #define THREAD_PER_BLOCK 32
3
4 __global__ void addition(int * a, int * b, int * c, unsigned int limit_d)
5 {
6     int id = blockIdx.x * THREAD_PER_BLOCK + threadIdx.x;
7     if(id < limit_d)
8         c[id] = a[id] + b[id];
9 }
10
11 int main(){
12     ...
13
14     // allocation de la mémoire
15     cudaMalloc((void**) &a_d, size * sizeof(int));
16     cudaMalloc((void**) &b_d, size * sizeof(int));
17     cudaMalloc((void**) &c_d, size * sizeof(int));
18
19     // transfert des données vers le processeur graphique
20     cudaMemcpy(a_d, a_h, size * sizeof(int), cudaMemcpyHostToDevice);
21     cudaMemcpy(b_d, b_h, size * sizeof(int), cudaMemcpyHostToDevice);
22
23     // appel du kernel
24     unsigned int blocks = (size) / THREAD_PER_BLOCK +
25         ((size) % THREAD_PER_BLOCK > 0);
26     addition<<<blocks, THREAD_PER_BLOCK>>>(a_d, b_d, c_d, size);
27
28     ...
29 }

```

L'allocation de la mémoire, le transfert des données et l'appel d'un kernel sont illustrés au listing 3.1. L'appel du kernel est effectué à la ligne 26 et se présente sous la forme :

```
kernel<<<Nombre_de_blocs,Nombre_de_threads_par_bloc>>>(arguments);
```

Pour mieux coller la représentation des données, il est possible de définir le nombre de blocs et le nombre de threads par bloc sous une représentation en 2D ou 3D. Il existe un type `dim3` qui permet cette représentation :

```
dim3 var(x,y,z)
```

En fait, lors de l'appel du kernel, les valeurs `Nombre_de_blocs` et `Nombre_de_threads_par_bloc` sont stockées dans deux variables de type `dim3`. Ainsi, l'exemple présent au listing 3.1 est un raccourci pour l'appel présent au listing 3.2.

Listing 3.2: Appel d'un kernel en utilisant des variables de type `dim3`

```

1
2     ...
3
4     // appel du kernel
5     dim3 blocks_3D(blocks,1,1);
6     dim3 threads_3D(THREAD_PER_BLOCK,1,1);
7     addition<<<blocks_3D,threads_3D>>>(a_d,b_d,c_d,size);
8
9     ...
10 }
```

La figure 3.4 présente la répartition des threads avec deux variables `dim3`, initialisées `(2,2,1)`. Cette fonctionnalité permet de préparer une grille de blocs et threads dont la géométrie correspond à celle des données. Les `dim3` sont intéressantes lorsque les données sont organisées en structure possédant plusieurs dimensions (matrice, cube ...). Ceci facilite le calcul de l'index d'accès aux données dans chacun des threads.

5. Le kernel:

En-tête de la fonction

Il existe trois types d'appel de kernel, chacun tant spécifié par un préfixe en tête de la définition du kernel :

- `__global__` : indique que le kernel est appelé depuis le CPU et exécuté sur le GPU
- `__device__` : indique que le kernel est appelé depuis le GPU et exécuté par le GPU
- `__host__` : indique que le kernel est appelé depuis le CPU et exécuté par le CPU (correspond aux procédures traditionnelles d'un CPU)

Il y a des limitations par rapport à la programmation sur CPU : un kernel ne peut être récursif et il ne peut retourner aucune valeur. Un exemple est présent au listing 3.1, la ligne 4.

Identification du thread

Les cœurs du GPU effectuent tous le même code et l'identifiant des threads est le seul élément qui permet de les différencier. Pour identifier le thread au sein du kernel, CUDA fournit trois données de type `dim3` qui identifient un bloc au sein de la grille et un thread au sein de son bloc :

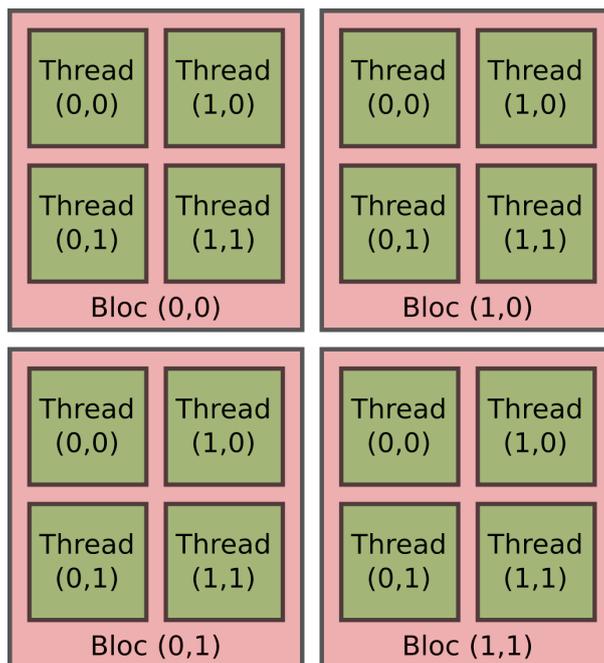


Figure 3.4: Rsltat de l'appel `kernel<<(2,2,1), (2,2,1)>>(arguments)` ;

- `threadId` : fournit l'index du thread dans le bloc.
- `blockId` : fournit l'index du bloc dans la grille.
- `blockDim` : fournit le nombre de threads par blocs.

Dans le cas d'une matrice (image 2D ...), il suffit d'utiliser ce code pour localiser l'lement associ au thread :

```
int Col = blockDim.x * blockIdx.x + threadIdx.x;
int Row = blockDim.y * blockIdx.y + threadIdx.y;
```

Dans le cas d'une image 3D, il est intressant d'utiliser les trois dimensions :

```
int Col = blockDim.x * blockIdx.x + threadIdx.x;
int Row = blockDim.y * blockIdx.y + threadIdx.y;
int Depth = blockDim.z * blockIdx.z + threadIdx.z;
```

La gestion de la mmoire

Comme cela a t prsent la section 3.3, le processeur graphique est pourvu de diffrents niveaux de mmoires, plus ou moins proches des cœurs, qui permettent de mieux grer la latence d'accs aux donnes. Certaines de ces mmoires sont alloues mme le kernel :

- La mmoire partage : cette mmoire permet de partager des donnes entre tous les threads d'un mme bloc. Pour dclarer une variable dans la mmoire partage, il suffit de faire prceder sa dclaration au sein du kernel par le mot-cl `__shared__`. La fonction `__syncthreads()` permet de synchroniser les threads d'un bloc et de s'assurer que les blocs de mmoire partage sont mis jour par tous les threads. Il y a trois cas d'accs optimiss la mmoire partage :

- Le plus rapide : tous les threads accèdent en lecture la même donnée. Un système de broadcast est mis en œuvre par le matériel.
- Des threads consécutifs accèdent des données consécutives : le décodage de l'adresse de la donnée n'est réalisé qu'une seule fois, les adresses suivantes sont incrémentées par le matériel.
- Le cas général : les threads accèdent aux données de façon aléatoire et chaque accès à une donnée nécessite un décodage spécifique de son adresse.

La durée de vie des données en mémoire partagée est celle du bloc.

- Les registres : cette mémoire est propre à un thread et possède sa propre durée de vie. Toute variable déclarée sans préfixe dans un kernel est prioritairement placée dans un registre.

Il existe d'autres mémoires possédant une latence plus faible que la mémoire globale, mais devant être initialisées avant l'appel au kernel :

- La mémoire constante : cette mémoire est accessible par tous les threads ainsi que par l'hôte. Elle possède la durée de vie de l'application. Une variable est dans la mémoire constante lorsque sa déclaration est précédée du mot-clé `__constant__`. Une telle variable ne peut être déclarée dans le corps d'une fonction. Cette mémoire agit comme la mémoire partagée pour ce qui en est de son accès : il y a lecture sans conflit de son contenu si tous les threads accèdent à un élément unique (broadcast), et il y a sérialisation des accès si des éléments différents sont accédés simultanément par différents threads. Elle est surtout utilisée pour les scalaires dont la valeur ne changera pas en cours d'exécution. Le transfert d'une donnée dans une telle mémoire nécessite l'appel par le CPU à la fonction :


```
cudaMemcpyToSymbol(const char * symbol, const void * src, size_t count,
size_t offset, enum cudaMemcpyKind kind)
```

 Cette fonction copie les `count` octets présents à l'adresse indiquée par `src` vers `symbol + offset`. `symbol` représente soit une variable résidant en mémoire constante ou globale, soit la chaîne de caractères identifiant la variable déclarée en mémoire globale ou constante. `kind` peut prendre comme valeur `cudaMemcpyHostToDevice` ou `cudaMemcpyDeviceToDevice`.
- La mémoire texture : c'est une mémoire optimisée pour un espace à deux dimensions. Les threads d'un même warp lisant des adresses proches auront des performances optimales. La lecture des mémoires par le mécanisme des textures peut être une alternative avantageuse à la lecture depuis les mémoires globale et constante. L'initialisation de la texture se fait à partir de l'hôte par la fonction :


```
cudaError_t cudaBindTextureToArray(const struct texture<T,dim,readMode> &
dst, const struct cudaArray * src, const struct cudaChannelFormatDesc & desc)
```

 Cette fonction permet de charger la mémoire texture `dst` sur le device par le tableau de deux dimensions `src` se trouvant en mémoire hôte. `desc` décrit la manière de lire la mémoire texture (cf. documentation ou M. Sidi). `dst` est une variable du type prédéfini `texture` :


```
struct texture<type, dim, cudaReadModeElementType> dst
```

 Tel que `type` représente le type de données (entier, virgule flottante ...), `dim` représente le nombre de dimensions de la texture (généralement égale à 2) et `cudaReadModeElementType` est le mode de lecture. La lecture de la texture par le kernel se fait via la fonction :


```
tex2D(dst, x, y)
```

 Tel que `dst` représente la texture initialisée et déclarée en haut et `x` et `y` représentent les coordonnées des éléments à l'intérieur de la texture.

Synchronisation des threads

Il est possible de synchroniser les threads d'un bloc à l'aide de la fonction `__syncthreads()`. Il n'est actuellement pas possible de synchroniser l'ensemble des blocs de threads à l'intérieur d'un kernel, mais la version 4.0 de CUDA, qui sortira prochainement, permettra une telle synchronisation.

Branchements

Les processeurs graphiques sont construits de telle manière que tous les cœurs d'un SM effectuent la même opération à un moment donné. S'il y a un branchement (`if`, un `while` offrant un nombre d'itération différent selon le thread ...) au sein d'un warp, les branches divergentes sont exécutées les unes après les autres. Ainsi, si un `if` offre une divergence au sein d'un warp, il faut deux fois plus de temps pour exécuter le warp. Les branchements sont donc viter.

La fusion des accès à la mémoire globale

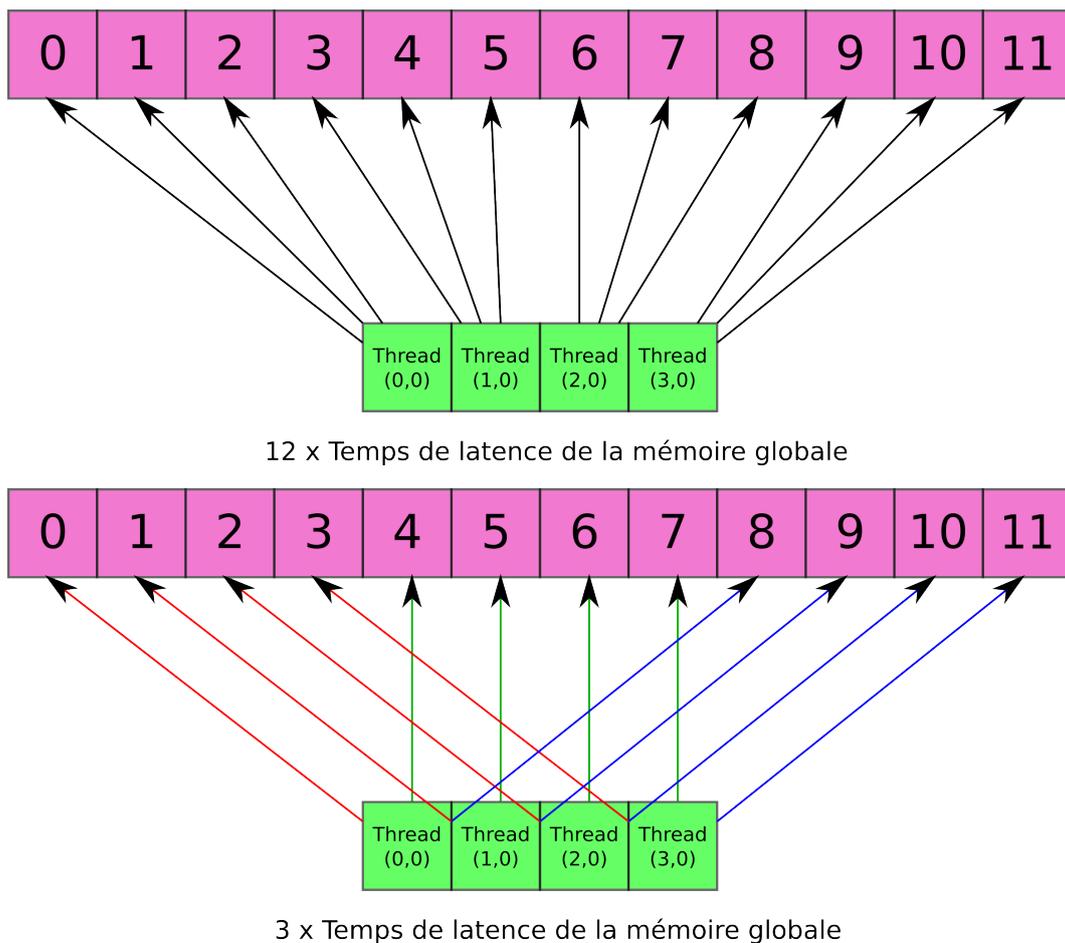


Figure 3.5: Présentation de l'accès fusionné à la mémoire globale : si des threads adjacents accèdent des espaces mémoire non-adjacents, les accès sont effectués les uns après les autres ; si les threads accèdent des espaces mémoires adjacents, les accès sont fusionnés et les données sont chargées en une fois.

L'accès à la mémoire globale implique une latence élevée de plus de 400 cycles d'horloge. Lorsque N threads adjacents d'un warp chargent N données adjacentes (cf. figure 3.5), les N accès sont fusionnés en un accès N éléments. Ceci permet de diviser la latence par quasiment N .

6. Transfert des données GPU \Rightarrow CPU:

La récupération des résultats sur le processeur central est assez simple, il suffit d'utiliser la nouvelle fonction `cudaMemcpy` présente au point 3 en utilisant la constante `cudaMemcpyDeviceToHost`.

7. Libération des ressources:

Enfin, toute ressource allouée doit être libre. Pour ce faire, CUDA fournit une version GPU de la fonction `free`¹⁰ :

```
cudaError_t cudaFree((void * devPtr))
```

`devPtr` indique l'espace mémoire à libérer.

Un exemple de transfert de données vers le CPU et de libération d'espace mémoire est présent au listing 3.3.

Listing 3.3: Addition de vecteurs : récupération du résultat et libération des ressources

```
1 int main(){
2     ...
3
4     //transfert des données vers le processeur central
5     cudaMemcpy(c_h, c_d, size * sizeof(int), cudaMemcpyDeviceToHost);
6
7     // libération des espaces mémoire alloués
8     cudaFree(a_d);
9     cudaFree(b_d);
10    cudaFree(c_d);
11
12    ...
13 }
```

8. Exemple complet

Listing 3.4: Addition de vecteurs : récupération du résultat et libération des ressources

```
1 #include <sys/time.h>
2 #include <stdio.h>
3
4 #define THREAD_PER_BLOCK 512
5
6 __global__ void addition(int * a, int * b, int * c, unsigned int limit_d)
7 {
8     // récupération de l'identifiant du thread
9     int id = blockIdx.x * THREAD_PER_BLOCK + threadIdx.x;
10    //si la taille n'est pas un multiple de la taille des warp,
11    //on vérifie que dans les limites des tableaux
12    if(id < limit_d){
13        //addition avec les accès aux données fusionnés
14        c[id] = a[id] + b[id];
15    }
16 }
```

¹⁰La fonction `free` permet de libérer une ressource allouée sur un processeur central.

```

17
18
19 int main(int argc, char * argv[]){
20     //allocation des ressources sur le processeur central
21     unsigned int size = atoi(argv[1]), i;
22     int * a_h = (int *) malloc(size * sizeof(int)), * a_d;
23     int * b_h = (int *) malloc(size * sizeof(int)), * b_d;
24     int * c_h = (int *) malloc(size * sizeof(int)), * c_d;
25
26     //allocation des ressources sur le processeur graphique
27     cudaMalloc((void**) &a_d, size * sizeof(int));
28     cudaMalloc((void**) &b_d, size * sizeof(int));
29     cudaMalloc((void**) &c_d, size * sizeof(int));
30     srand(1234);
31     for(i=0;i<size;++i){
32         a_h[i] = rand()%size;
33         b_h[i] = rand()%size;
34     }
35
36     //allocation des ressources pour mesurer le temps
37     unsigned long long time = 0;
38     struct timeval tv1,tv2;
39     gettimeofday(&tv1, NULL);
40
41     //transfert des données sur le GPU
42     cudaMemcpy(a_d, a_h, size * sizeof(int), cudaMemcpyHostToDevice);
43     cudaMemcpy(b_d, b_h, size * sizeof(int), cudaMemcpyHostToDevice);
44
45     //appel au kernel
46     unsigned int blocks = (size) / THREAD_PER_BLOCK +
47         ((size) % THREAD_PER_BLOCK > 0);
48     addition<<<blocks,THREAD_PER_BLOCK>>>(a_d,b_d,c_d,size);
49
50     //transfert des résultats sur le CPU
51     cudaMemcpy(c_h, c_d, size * sizeof(int), cudaMemcpyDeviceToHost);
52
53     //calcul du temps
54     gettimeofday(&tv2, NULL);
55     time+=(tv2.tv_sec-tv1.tv_sec) * 1000000L + (tv2.tv_usec-tv1.tv_usec);
56     printf("Temps de calcul GPU : %llu usec \n", time);
57
58     //libération des ressources sur le GPU
59     cudaFree(a_d);
60     cudaFree(b_d);
61     cudaFree(c_d);
62
63     //libération des ressources sur le CPU
64     free(a_h);
65     free(b_h);
66     free(c_h);
67 }

```

Pour ceux qui veulent aller plus loin avec CUDA, ils peuvent consulter quelques exemples d'utilisation de GPU via CUDA pour le traitement d'images et de vidéos dans [10],[11] et [12].

3.4.3 OpenCL

OpenCL (Open Computing Language) se veut au GPGPU ce qu'OpenGL est au rendu 3D, savoir un standard ouvert. OpenCL se compose d'une API et d'un langage de programmation driv du C. Son but est de faciliter la programmation d'applications utilisant les divers processeurs

disponibles dans les machines : les CPU, les GPU, et peut être tendu vers d'autres types de processeurs.

Le projet de création d'OpenCL a été initié par Apple, qui en a ensuite confié la gestion au Khronos Group [13]. Le projet a depuis été rejoint par les principaux acteurs du domaine : AMD, nVidia et Intel. Leurs futurs produits devraient être compatibles avec les spécifications du standard. On peut donc s'attendre à ce qu'OpenCL devienne un standard très utilisé. La première version est disponible depuis Avril 2010. OpenCL bénéficie essentiellement de son indépendance vis-à-vis des API 3D, sa compatibilité avec les cartes de tous les constructeurs ainsi que sa portabilité sur différentes plateformes.

Pour ceux qui veulent aller plus loin avec OpenCL, ils peuvent consulter quelques exemples d'utilisation de GPU via CUDA pour le traitement d'images médicales dans [14]

3.5 Exploitation des architectures multi-cœurs hétérogènes

Les processeurs graphiques ont permis de fournir une solution très efficace pour l'accélération des applications intensives en calcul. Cependant, cette solution peut être encore améliorée par l'exploitation simultanée des cœurs CPU et GPU multiples. Dans ce contexte, il existe différents travaux pour l'exploitation des plateformes multi-cœurs et hétérogènes, tels que : StarPU, StarSS et GrandCentralDispatch.

3.5.1 StarPU

La bibliothèque StarPU [15], développée à l'INRIA Bordeaux (France), permet d'offrir un support exécutif unifié pour exploiter les architectures multi-cœurs hétérogènes, tout en s'affranchissant des difficultés liées à la gestion des transferts de données. StarPU propose par ailleurs plusieurs stratégies d'ordonnancement efficaces et offre la possibilité d'en concevoir aisément de nouvelles. StarPU se base pour son fonctionnement sur deux structures principales : la codelet et les tâches.

- **La codelet :** permet de préciser sur quelles architectures le noyau de calcul peut s'effectuer, ainsi que ses implémentations associées (CPU, GPU ...).
- **Les tâches :** elles consistent à appliquer la codelet sur l'ensemble des données, en utilisant la stratégie d'ordonnancement sélectionnée (StarPU fournit une interface qui permet le partitionnement des données).

L'évolution d'une tâche depuis sa soumission jusqu'à la notification de sa terminaison l'application peut être résumé en six étapes :

1. **Soumission :** StarPU reçoit toutes les tâches soumises à partir de l'application.
2. **Sélection des ressources de calcul :** l'ordonnanceur distribue les tâches selon la politique d'ordonnancement sélectionnée. Si la tâche a été attribuée au GPU : la description de la tâche est donc envoyée au pilote associé à ce GPU (chaque unité de calcul est associée à un pilote spécifique). Dès que disponible, le pilote réclame une nouvelle tâche à l'ordonnanceur.
3. **Chargement des données :** pour les tâches attribuées au GPU, les données doivent être transférées vers la carte graphique. StarPU permet de masquer ces transferts puisqu'il dispose d'une bibliothèque qui met en œuvre une mémoire partagée virtuelle (ou DSM pour Distributed Shared Memory). La programmation des mouvements de données est donc transparente pour le programmeur.
4. **Calcul StarPU :** une fois les données chargées, les tâches peuvent être exécutées sur les ressources hétérogènes de calcul indiquées dans la codelet.

5. **Rapatriement des rsultats** : une fois les calculs achevs, les rsultats peuvent tre rapatriés vers la mmoire hte.
6. **Libration des ressources** : une fois que la tche est terminée, StarPU en est notifié et lance une autre tche en fonction des dépendances des tches et de sa politique d'ordonnancement.

Un exemple d'utilisation de StarPU pour le traitement d'images est décrit dans [16],[17] et [18].

3.5.2 StarSs

StarSs [19] a été développé à l'université de Catalogne (Barcelone, Espagne). Il permet de fournir un modèle de programmation flexible pour les architectures multi-cœurs. Ainsi, il est composé de six composantes principales: CellSs, SPMSs, GPUSs, ClearSpeedSs, ClusterSs et GridSs.

1. **CellSs** : permet une exploitation automatique du parallélisme fonctionnel à partir d'un programme séquentiel à travers les éléments de traitement de l'architecture Cell BE.
2. **SPMSs** : permet d'exploiter de manière automatique et parallèle les cœurs multiples des processeurs SMP (Symmetric multiprocessing) à partir d'un code séquentiel écrit par le programmeur.
3. **ClearSpeedSs** : permet aussi une exploitation automatique et parallèle des cœurs multiples des architectures ClearSpeed.
4. **GPUSs** : permet l'exploitation des processeurs graphiques (GPU) multiples.
5. **ClusterSs** : permet l'exploitation parallèle des serveurs présents dans les clusters.
6. **GridSs** : permet une exploitation parallèle des ressources de calcul disponibles dans une grille.

3.5.3 GrandCentralDispatch

Grand Central Dispatch [20] est une technologie développée par Apple pour mieux exploiter les processeurs multi-cœurs dans les plateformes Mac. Cette nouvelle architecture permet aux développeurs d'utiliser pleinement la puissance de calcul des processeurs multi-cœurs. Elle travaille en distribuant de manière efficace les traitements (processus) aux différents cœurs. Le code de Grand Central Dispatch est ouvert depuis septembre 2009.

Une comparaison entre ces technologies d'exploitation de plateformes multi-cœurs (StarPU, StarSs et GrandCentralDispatch) nous permet de distinguer que StarPU est la plus prometteuse grâce à son exploitation de l'intégralité des ressources hétérogènes de calcul (multi-CPU/multi-GPU), l'inverse de StarSs et GrandCentralDispatch qui ne permettent qu'une exploitation distincte des cœurs GPUs multiples ou CPU multiples.

Bibliography

- [1] OpenMP ARB. The OpenMP API specification for parallel programming. <http://openmp.org/wp/>.
- [2] GCC team. Welcome to the home of GOMP. <http://gcc.gnu.org/projects/gomp/>.
- [3] Dartmouth College. What is OpenMP ? http://www.dartmouth.edu/~rc/classes/intro_openmp/.
- [4] Lawrence Livermore National Laboratory Blaise Barney. OpenMP. <https://computing.llnl.gov/tutorials/openMP/>.
- [5] DeinoMPI. The Great and Terrible implementation of MPI-2. http://mpi.deino.net/mpi_functions/.
- [6] Argonne National Laboratory. Web pages for MPI and MPE. <http://www.mcs.anl.gov/research/projects/mpi/www/>.
- [7] Open MPI. Open MPI v1.5.4 documentation. <http://www.open-mpi.org/doc/v1.5/>.
- [8] OpenGL. OpenGL Architecture Review Board: ARB vertex program. Revision 45. <http://oss.sgi.com/projects/ogl-sample/registry/>, 2004.
- [9] SGI. Silicon Graphics. OpenGL-the industry’s foundation for high performance graphics. <http://www.sgi.com/products/software/opengl/>, 1992.
- [10] P. d. C. Possa, S. A. Mahmoudi, N. Harb, and C. Valderrama. A new self-adapting architecture for feature detection. In *22nd International Conference on Field Programmable Logic and Applications (FPL)*, pages 643–646, Aug 2012.
- [11] S. A. Mahmoudi and P. Manneback. Multi-gpu based event detection and localization using high definition videos. In *2014 International Conference on Multimedia Computing and Systems (ICMCS)*, pages 81–86, April 2014.
- [12] Sidi Ahmed Mahmoudi, Michal Kierzynka, and Pierre Manneback. *Real-Time GPU-Based Motion Detection and Tracking Using Full HD Videos*, pages 12–21. Springer International Publishing, Cham, 2013.
- [13] KHRONOS GROUP. Khronos group : Open standards for media authoring and acceleration. <http://www.khronos.org/about/>, 2000.
- [14] Mohamed Amine Larhmam, Sidi Ahmed Mahmoudi, Mohammed Benjelloun, Saïd Mahmoudi, and Pierre Manneback. *A Portable Multi-CPU/Multi-GPU Based Vertebra Localization in Sagittal MR Images*, pages 209–218. Springer International Publishing, Cham, 2014.

- [15] Cdric Augonnet, Samuel Thibault, Raymond Namyst, and Pierre-Andr Wacrenier. StarPU: A Unified Platform for Task Scheduling on Heterogeneous Multicore Architectures. *In Concurrency and Computation: Practice and Experience, Euro-Par 2009, best papers issue*, pages 863–874, 2009.
- [16] S.A. Mahmoudi, P. Manneback, C. Augonnet, S. Thibault, et al. Traitements d’images sur architectures parallèles et hétérogènes. *Technique et Science Informatiques, Revue des sciences et technologies de l’information*, 2012.
- [17] S.A. Mahmoudi, P. Manneback, C. Augonnet, and S. Thibault. Détection optimale des coins et contours dans des bases d’images volumineuses sur architectures multicœurs hétérogènes. *In 20ème Rencontres francophones du parallélisme (RenPar’20)*, 2011.
- [18] Sidi Ahmed Mahmoudi and Pierre Manneback. *Multi-CPU/Multi-GPU Based Framework for Multimedia Processing*, pages 54–65. Springer International Publishing, Cham, 2015.
- [19] Eduard Ayguad, Rosa M Badia, Francisco D Igual, Jesus Labarta, Rafael Mayo, and Enrique S Quintana-Orti. An Extension of the StarSs Programming Model for Platforms with Multiple GPUs. *Proceedings of the 15th International Euro-Par Conference on Parallel Processing. Euro-Par’09*, pages 851–862, 2009.
- [20] Apple. Grand Central Dispatch. A better way to do multicore, 2009.