

# Techniques d'optimisation de la parallélisation

## Examen écrit

Master 1 MIHPS

Nom : \_\_\_\_\_

Prénom : \_\_\_\_\_

### I – QCM

Ce QCM comporte en tout 20 questions. Il n'y a pas de pénalité associée aux mauvaises réponses. **Si une question accepte plusieurs réponses, elle le mentionne explicitement en précisant leur nombre.** Dans les questions suivantes **seq(n)** correspond à la partie séquentielle, **par(n)** à la partie parallèle et enfin **comm(n,p)** aux communications **n**'étant la taille du problème et **p** le nombre de cœurs. De plus la notation **S(n,p)** est associée au speedup.

Question n°1 :	Les principaux facteurs limitant le parallélisme d'une application sont ( <b>deux réponses attendues</b> )		
<input type="checkbox"/> Les communications	<input type="checkbox"/> Le compilateur	<input type="checkbox"/> Les threads	
<input type="checkbox"/> La météo	<input type="checkbox"/> La fraction parallèle	<input type="checkbox"/> La fraction séquentielle	
Question n°2 :	Deux boucles <i>for</i> imbriquées ont une complexité en ...		
<input type="checkbox"/> $\Theta(n \ln_2(n))$	<input type="checkbox"/> $\Theta(n)$	<input type="checkbox"/> $\Theta(n^2)$	
<input type="checkbox"/> $\Theta(1)$	<input type="checkbox"/> $\Theta(\sqrt{n})$	<input type="checkbox"/> $\Theta(n^{\log(n)})$	
Question n°3 :	Une recherche dans une liste doublement chaînée à une complexité en ...		
<input type="checkbox"/> $\Theta(n \ln_2(n))$	<input type="checkbox"/> $\Theta(n)$	<input type="checkbox"/> $\Theta(n^2)$	
<input type="checkbox"/> $\Theta(1)$	<input type="checkbox"/> $\Theta(\sqrt{n})$	<input type="checkbox"/> $\Theta(n^{\log(n)})$	
Question n°4 :	En MPI les communications collectives ont une complexité en ...		
<input type="checkbox"/> $\Theta(1)$	<input type="checkbox"/> $\Theta(n)$	<input type="checkbox"/> $\Theta(\sqrt{n})$	
<input type="checkbox"/> $\Theta(n \ln(n))$	<input type="checkbox"/> $\Theta(e^n)$	<input type="checkbox"/> $\Theta(n^2)$	
Question n°5 :	Si le déroulement d'un algorithme peut se représenter sous la forme d'un arbre binaire alors sa complexité est probablement proche de ...		
<input type="checkbox"/> $\Theta(n \ln_2(n))$	<input type="checkbox"/> $\Theta(n^2)$	<input type="checkbox"/> $\Theta(n)$	
<input type="checkbox"/> $\Theta(1)$	<input type="checkbox"/> $\Theta(\sqrt{n})$	<input type="checkbox"/> $\Theta(n^{\log(n)})$	
Question n°6 :	$S(n, p) \leq \frac{1}{s + \frac{(1-s)}{p}} \quad \text{avec} \quad s = \frac{seq(n)}{seq(n) + par(n)} \quad \text{c'est ...}$		
<input type="checkbox"/> L'efficacité	<input type="checkbox"/> La loi de Gustafson-Barsis	<input type="checkbox"/> La loi de Amdahl	
<input type="checkbox"/> La métrique de Karp-Flatt	<input type="checkbox"/> La formule du speedup	<input type="checkbox"/> L'équation de Maxwell	

Question n°7 :	L'efficacité d'un programme se calcule en faisant :		
<input type="checkbox"/> $\frac{1}{S(n, p)}$	<input type="checkbox"/> $\frac{p}{S(n, p)}$	<input type="checkbox"/> $\sqrt{\frac{seq(n)}{par(n)}}$	
<input type="checkbox"/> $\frac{seq(n)}{p}$	<input type="checkbox"/> $\frac{e^{seq(n) + comm(n, p)}}{p}$	<input type="checkbox"/> $\frac{S(n, p)}{p}$	
Question n°8 :	En vous servant de la formule de la question 6 quel est le speedup maximal que l'on peut espérer si 20 % du temps d'exécution est séquentiel ?		
<input type="checkbox"/> 0,25	<input type="checkbox"/> 0,5	<input type="checkbox"/> 2	
<input type="checkbox"/> 4	<input type="checkbox"/> 5	<input type="checkbox"/> 20	
Question n°9 :	D'après la loi d'Amdahl, le speedup d'un programme est limité par (deux réponses attendues) ...		
<input type="checkbox"/> La vitesse de la lumière	<input type="checkbox"/> La fraction séquentielle	<input type="checkbox"/> La fraction parallèle	
<input type="checkbox"/> Les communications	<input type="checkbox"/> $\frac{1}{s}$	<input type="checkbox"/> Le programmeur	
Question n°10 :	Un programme passe à l'échelle d'autant plus que ....		
<input type="checkbox"/> La machine est grande	<input type="checkbox"/> Son problème est petit	<input type="checkbox"/> Son problème est grand	
<input type="checkbox"/> La fraction séquentielle est importante	<input type="checkbox"/> Le réseau de communication est lent	<input type="checkbox"/> Les communications sont prépondérantes	
Question n°11 :	Du point de vue optimisation un programme parallèle est idéal si ... (deux réponses attendues)		
<input type="checkbox"/> Sa complexité est exponentielle	<input type="checkbox"/> La fraction séquentielle est importante	<input type="checkbox"/> Il n'y a pas de communications	
<input type="checkbox"/> Le problème est facilement divisible	<input type="checkbox"/> Sa scalabilité est limitée	<input type="checkbox"/> Les communications sont prépondérantes	
Question n°12 :	Quels éléments parmi les suivants peuvent nuire à une mesure (quatre réponses attendues)		
<input type="checkbox"/> Utiliser un compteur précis	<input type="checkbox"/> La répéter de nombreuses fois	<input type="checkbox"/> La faire une seule fois	
<input type="checkbox"/> Utiliser un compteur peu précis	<input type="checkbox"/> Mesurer un temps trop court	<input type="checkbox"/> Les autres programmes s'exécutant sur la machine	
Question n°13 :	La scalabilité faible se calcule pour une taille de problème ...		
<input type="checkbox"/> variable	<input type="checkbox"/> très petite	<input type="checkbox"/> constante	

<input type="checkbox"/> très grande	<input type="checkbox"/> proportionnelle au nombre de cœurs	<input type="checkbox"/> inversement proportionnelle
--------------------------------------	---	--

<b>Question n°14 :</b>	Lequel de ces programmes permet d'obtenir un profil des communications MPI ?	
<input type="checkbox"/> gdb	<input type="checkbox"/> mpich	<input type="checkbox"/> Scalasca
<input type="checkbox"/> gprof	<input type="checkbox"/> gcc	<input type="checkbox"/> valgrind
<b>Question n°15 :</b>	Lors d'une décomposition de domaine en 2X2 cartésien avec des conditions de bord « nulles » combien de communications MPI il y a t'il au total par phase d'échange ? Chaque nœud échange avec ses 8 voisins.	
<input type="checkbox"/> 4	<input type="checkbox"/> 8	<input type="checkbox"/> 12
<input type="checkbox"/> 16	<input type="checkbox"/> 18	<input type="checkbox"/> 32
<b>Question n°16 :</b>	Quel est l'ordre de grandeur de la latence d'une communication <u>intra-nœud</u> ?	
<input type="checkbox"/> 2 ns	<input type="checkbox"/> 500 ns	<input type="checkbox"/> 2 μs
<input type="checkbox"/> 500 μs	<input type="checkbox"/> 2 ms	<input type="checkbox"/> 500 ms
<b>Question n°17 :</b>	L'analyse d'un programme MPI met en évidence trois MPI_Bcast (voir annexes pour manuel de MPI_Bcast) lors de la phase d'initialisation. Que dois-je faire pour optimiser ce programme ?	
<input type="checkbox"/> Les supprimer	<input type="checkbox"/> les remplacer par des MPI_Isend et MPI_Irecv	<input type="checkbox"/> les remplacer par des MPI_Send et MPI_Recv
<input type="checkbox"/> Fuir	<input type="checkbox"/> Déplacer l'initialisation à la fin	<input type="checkbox"/> Rien, ils ne sont pas un problème,
<b>Question n°18 :</b>	Un programme MPI s'exécute en 25 secondes sur 8 cœurs et pour un même problème en 100 secondes sur 1 cœur, L'efficacité de ce programme sur huit cœurs et de :	
<input type="checkbox"/> 0,1	<input type="checkbox"/> 0,125	<input type="checkbox"/> 0,25
<input type="checkbox"/> 0,5	<input type="checkbox"/> 1	<input type="checkbox"/> 2
<b>Question n°19 :</b>	On dit qu'il y a false-sharing entre deux variables lorsqu'elles sont accédées simultanément par deux threads et qu' ...	
<input type="checkbox"/> elles sont sur la même ligne de cache	<input type="checkbox"/> elles sont sur deux lignes de caches voisines	<input type="checkbox"/> elles sont dans le même tableau

<input type="checkbox"/> elles ont la même adresse	<input type="checkbox"/> elles ont la même valeur	<input type="checkbox"/> elles ont même nom
<b>Question n°20 :</b>	Un processeur hyper-threadé ...	
<input type="checkbox"/> peut exécuter deux threads sur un même cœur physique	<input type="checkbox"/> a physiquement deux fois plus de cœurs	<input type="checkbox"/> a deux fois plus d'unités flottantes.
<input type="checkbox"/> est deux fois plus rapide	<input type="checkbox"/> peut changer de fréquence	<input type="checkbox"/> fait pour OpenMP

## II – Questions

1) Quelle est la différence entre un spinlock et un mutex ? Dans quel cas de figure les utilisez-vous ?

2) Donnez un exemple de false-sharing.

3) Implémentez une réduction MPI en utilisant un arbre binaire (code C)

*values* : Valeur à réduire en somme.

*size* : nombre de valeurs.

*comm* : communicateur sur lequel effectuer la réduction.

*root* : rang du processus racine.

```
void my_reduce( int *values, int size, int root, MPI_Comm comm )
```

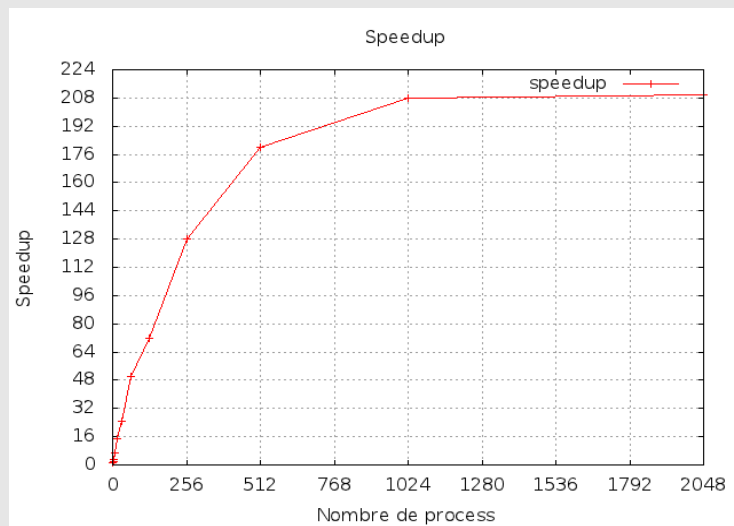
```
{
```

```
}
```

4) Définissez ce qu'est une architecture NUMA, quelles précautions doit-on prendre lorsqu'on les utilise en mode multi-threads ?

5) Quelle est la différence entre OpenMP et MPI ? Quel est l'intérêt pour les architectures modernes ?

5)



Répondez aux questions suivantes en utilisant cette courbe de speedup.

Quelle est l'efficacité de ce code sur 1024 cœurs ?

Comment pourriez-vous estimer la fraction séquentielle de ce programme ?

Que répondriez-vous au propriétaire de ce code vous expliquant que son code va plus vite sur 2048 cœurs que sur 512 ? Quelle métrique de performance utiliserez-vous donnez un exemple argumenté.



### III – Critiques de codes

Dans cette partie vous allez devoir commenter des extraits de code en décrivant ce qui les rend peu performants et en proposant une correction améliorant leurs performances. Par souci de simplicité seules les parties intéressantes sont retenues.

Les manuels des fonctions utiles sont à la fin du sujet.

**Vous ne devez pas écrire de code, une simple description de la correction suffit.**

#### **Code n°1 :**

```
if( rank == 0 )
{
    for( i = 0 ; i < buffer_size ; i++ )
        MPI_Send( buffer + i, 1, MPI_FLOAT, 1, 0,
MPI_COMM_WORLD);
}

if( rank == 1 )
{
    MPI_Status status;
    for( i = 0 ; i < buffer_size ; i++ )
        MPI_Recv( buffer + i, 1, MPI_FLOAT, 0, 0,
MPI_COMM_WORLD, &status);
}
```

Quels sont les défauts de ce code ? De quelle manière auriez-vous procédé ?

## **Code n°2 :**

```
MPI_Comm_rank( MPI_COMM_WORLD, &rank);
MPI_Comm_size( MPI_COMM_WORLD, &size);

if( rank == 0 )
{
    for( i = 0 ; i < size, i++ )
        MPI_Send( &resultat_important, 1, MPI_FLOAT, i, 0,
MPI_COMM_WORLD);
}
else
{
    MPI_Status status;
    MPI_Recv( &resultat_important, 1, MPI_FLOAT, 0, 0,
MPI_COMM_WORLD, &status);
}
```

Quels sont les défauts de ce code ? De quelle manière auriez-vous procédé ?

### **Code n°3 :**

```
MPI_Comm_rank( MPI_COMM_WORLD, &rank);
MPI_Comm_size( MPI_COMM_WORLD, &size);

if( rank == 0 )
{
    resultat_partiel = resultat_local;
    MPI_Send( &resultat_partiel, 1, MPI_FLOAT, (rank + 1) %
size, 0, MPI_COMM_WORLD);
}
else
{
    MPI_Status status;
    MPI_Recv( &resultat_partiel, 1, MPI_FLOAT, rank-1, 0,
MPI_COMM_WORLD, &status);
    resultat_partiel = resultat_partiel + resultat_local;
    MPI_Send( &resultat_partiel, 1, MPI_FLOAT, (rank + 1) %
size, 0, MPI_COMM_WORLD);
}

if( rank == 0 )
{
    MPI_Status status;
    MPI_Recv( &resultat_total, 1, MPI_FLOAT, size - 1, 0,
MPI_COMM_WORLD, &status);
}
```

Quels sont les défauts de ce code ? De quelle manière auriez-vous procédé ?

#### **Code n°4 :**

```
#include <pthread.h>
```

```
void * calcul_thread( void *a )
{
    pthread_mutex_lock( &verrou_travaux );
    struct travail_thread *travail =
        recuperer_travail( &liste_travaux );
    effectuer_travail( travail );
    pthread_mutex_unlock( &verrou_travaux );
    return NULL;
}

int main( int argc, char **argv )
{
    pthread_t threads[10];

    int i = 0;

    for( i = 0 ; i < 10 ; i++ )
        pthread_create( &threads[i], NULL, calcul_thread, NULL );
    for( i = 0 ; i < 10 ; i++ )
        pthread_join( threads[i], NULL );

    return 0;
}
```

Quels sont les défauts de ce code ? De quelle manière auriez-vous procédé ?

# MPI\_Send

Performs a basic send

## Synopsis

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

## Input Parameters

<b>buf</b>	initial address of send buffer (choice)
<b>count</b>	number of elements in send buffer (nonnegative integer)
<b>datatype</b>	datatype of each send buffer element (handle)
<b>dest</b>	rank of destination (integer)
<b>tag</b>	message tag (integer)
<b>comm</b>	communicator (handle)

## Notes

This routine may block until the message is received.

# MPI\_Recv

Basic receive

## Synopsis

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

## Output Parameters

<b>buf</b>	initial address of receive buffer (choice)
<b>status</b>	status object (Status)

## Input Parameters

<b>count</b>	maximum number of elements in receive buffer (integer)
<b>datatype</b>	datatype of each receive buffer element (handle)
<b>source</b>	rank of source (integer)
<b>tag</b>	message tag (integer)
<b>comm</b>	communicator (handle)

## MPI\_Reduce

Reduces values on all processes to a single value using given MPI operation.

### Synopsis

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

**Input**

### Parameters

**sendbuf** address of send buffer (choice)

**count** number of elements in send buffer (integer)

**datatype** data type of elements of send buffer (handle)

**op** reduce operation (handle)

**root** rank of root process (integer)

**comm** communicator (handle)

### Output Parameter

**recvbuf**

address of receive buffer (choice, significant only at `root`)

### Algorithm

This implementation currently uses a simple tree algorithm.

## MPI\_Bcast

Broadcasts a message from the process with rank "root" to all other processes of the group.

### Synopsis

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

### Input/output Parameters

**buffer** starting address of buffer (choice)

**count** number of entries in buffer (integer)

**datatype** data type of buffer (handle)

**root** rank of broadcast root (integer)

**comm** communicator (handle)

### Algorithm

If the underlying device does not take responsibility, this function uses a tree-like algorithm to broadcast the message to blocks of processes. A linear algorithm is then used to broadcast the message from the first process in a block to all other processes. `MPIR_BCAST_BLOCK_SIZE` determines the size of blocks. If this is set to 1, then this function is equivalent to using a pure tree algorithm. If it is set to the size of the group or greater, it is a pure linear algorithm. The value should be adjusted to determine the most efficient value on different machines.