

# Programming parallel computers

## Introduction

Xavier JUVIGNY, SN2A, DAAA, ONERA  
[xavier.juvigny@onera.fr](mailto:xavier.juvigny@onera.fr)

Course Parallel Programming  
- January 8th 2023 -

<sup>1</sup> ONERA, <sup>2</sup> DAAA

# Table of contents

---

# Overview

---

# Parallel architecture

---

## The main story

- Processors with multiples computing cores for faster computation ;
- Using simultaneously many cores for an unique application ;
- Performance benchmark in scientific computing is given by the number of FLoating Operations Per Seconds (FLOPS)

## Hardware implementation

- Many computing cores sharing a same main memory inside a computer ;
- Using many computers linked with a fast specialized ethernet connection ;
- Mixing both technologies above ;

# Interests of parallel architecture ?

---

- **Gordon Moore's "Law"** : In 1965, Gordon Moore (one of Intel's founder) observes that the number of transistors for each generation of processors double in heigten months, doubling the computing power ;
- In fact, **it isn't a law**, but it has been used by processors builder as a map road until 2000 years to raise the frequency of computing cores ;
- **Limitations of Gordon Moore Law** : The miniaturisation of transistors and the raising of their frequencies raises the heat inside the processor. Moreover, the miniaturisation is now at molecular scale and one must consider quantum effects (as tunnel effect) when making a processor ;
- Nowaday, **the Moore law is always verified**, but one don't double now the number of transistors inside a computing core but raise the number of computing core inside a processor or a computer.

# Parallel computing example (1)

## Control command

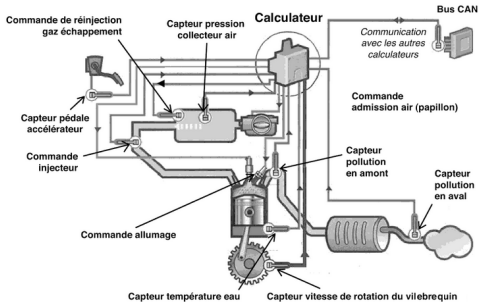


Figure – Car's control command

- Many tiny computers specialized for specific tasks : ABS, motor optimization, lighting, climatisation, wheel pressure optimization, mixing fuel/air, battery optimization, and so on.
- Computation must be terminated in constraint times ;
- Lot of parameters are interdependants (external air temperature, wheel pressure, optimal speed and oil consumption)

# Parallel computing example (1)

## Control command (continuation)



Figure – None Control Command of a reactor

- Another control command : managing nuclear power reactors ;
- High real time constraint algorithms ;
- Lot of complex computations ;
- One computing core isn't enough to satisfy tiny real time constraints ;
- **Solution** : Concurrency execution of interdependant tasks on many computing cores ;

# Parallel computation example (2)

## Physical Simulation

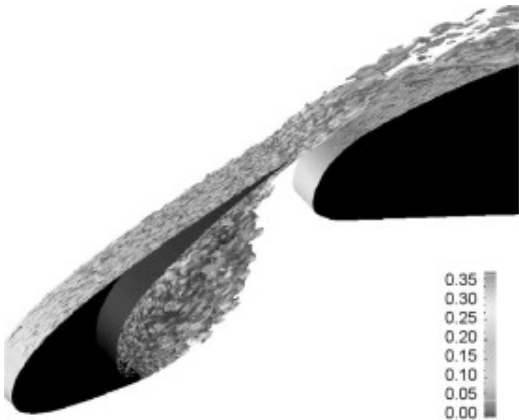


Figure – Turbulent noise generated by a plane's slate wing

- Turbulence : very small phenomena (millimeter scale). Need a mesh with lot of tiny triangles ;
- Typically, the mesh must contain five to ten billions of vertices with five unknown variables for each vertex ;
- Minimum memory requirement : 7 To
- Sequential computation time : 23 days to simulate  $\frac{1}{100}$  e seconds



# Parallel computation example (3)

## Artificial intelligence

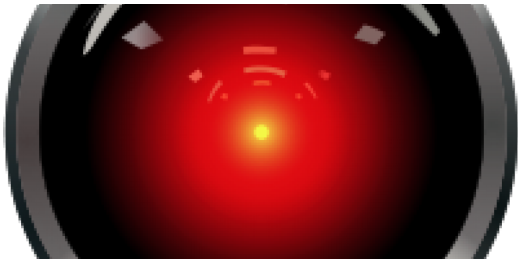


Figure – A very famous artificial intelligence (HAL)

- Deep learning used in AI to categorize

pictures, automatic translations, Cancerous cells detection, automatic vehicles, and so.

- In sequential required more than one year to learn ;
- with GPGPU, required about one month ;
- **March 2016** : Alphago wins versus world GO human champion (supervised learning) ;
- **October 2017** : Alphago zero wins versus alphago at 100 games to zero (deep learning).

# Parallel computation (4)

## Picture treatment

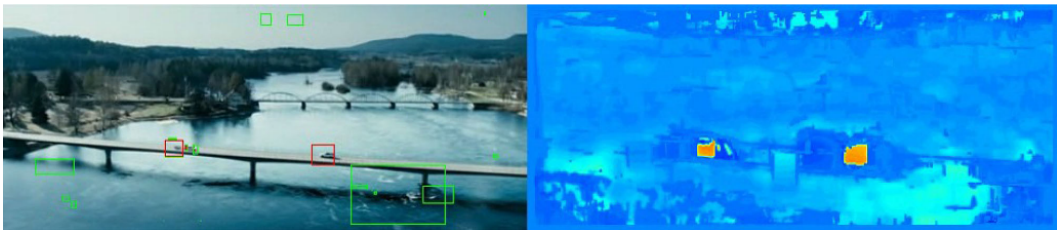


Figure – Real time constraint treatment of a video with 30 frames/s (resolution  $1920 \times 1080$  pixels)

- Needed for optical captors for navigation of autonomous vehicles, for super-resolution picture issued from low resolution video, and so on.
- Lot of computations needed (PDE equation to solve);
- Must use GPGPU and parallel algorithms to have real time constraint;

# Top 10 of supercomputers (June 2020)

Name	Core	Perf. (TFlops)	Constructor	Country	Power (kW)
<b>Fugaku</b>	7 299 072	415 530	Fujitsu	Japan	28 335
<b>Summit</b>	2 414 592	148 600	IBM	USA	10 096
<b>Sierra</b>	1 572 480	94 640	IBM/NVidia	USA	7 438
<b>Sunway TaihuLight</b>	10 649 600	93 014	NRCPC	China	15 371
<b>Tianhe-2A</b>	4 981 760	61 444	NUDT	China	18 482
<b>HPC5</b>	669 760	35 450	Dell EMC	Italy	2 252
<b>Selene</b>	277 760	27 580	Nvidia	USA	1 344
<b>Frontera</b>	448 448	23 516	Dell EMC	USA	?
<b>Marconi-100</b>	347 776	21 640	IBM	Italy	1 476
<b>Frontier</b>	591 872	1 102	HPE	USA	21 000

**Remarque** : Nowadays, we look for Flops/Watt performances

# Overview

---

# shared memory architecture

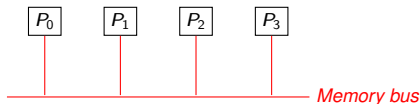


Figure – Simplified scheme of memory shared parallel architecture

Many computing cores shared the same main memory

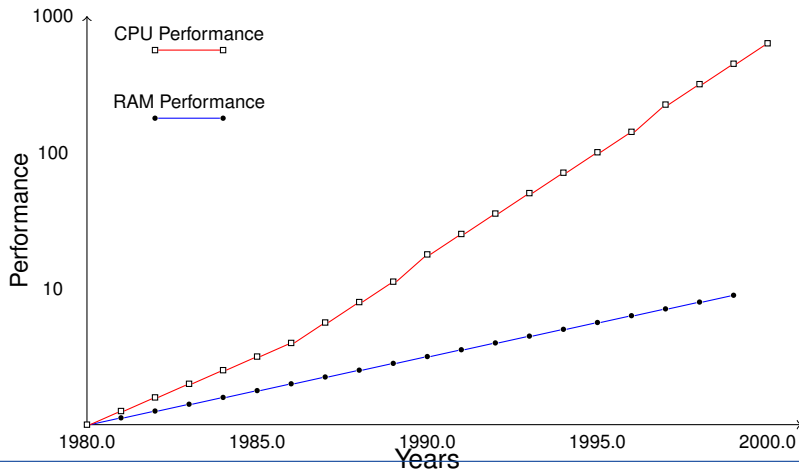
## Examples

- The recent multi-cores processors ;
- The graphic cards with 3D acceleration ;
- The phones, pad, etc.

## Memory access problem

- Optimization of memory access.
- Simultaneous read/write accesses at same memory location.

# Memory access



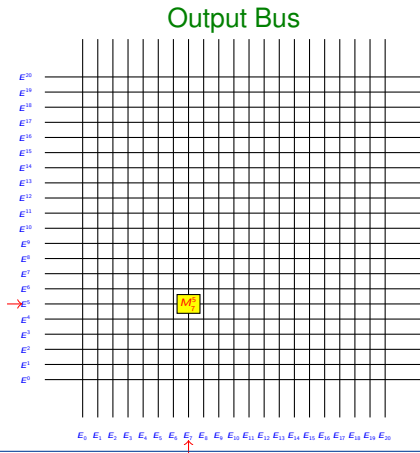
# Latency memory example (Haswell architecture)

Level	Size	Latency (cycles)	Physical location
<b>L1 Cache</b>	16/16 ko	4	In each core
<b>L2 Cache</b>	256 ko	12	Shared by two cores
<b>L3 Cache</b>	12 Mo	21	Shared by all cores
<b>Ram</b>	32 Go	117	SRAM on mother board
<b>Swap</b>	100+ Go	10 000	Hard disk or SSD

## Conclusion

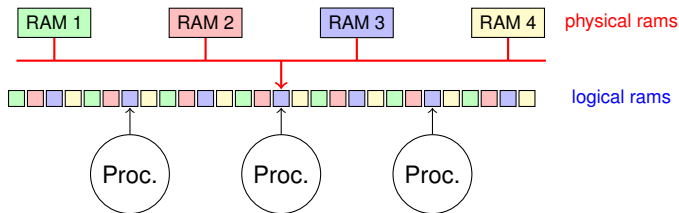
- Memory is more and more slower comparing to the instruction execution of the processor ;
- It's worst with multicore architecture !

# How works a RAM ?





# Interleaved rams

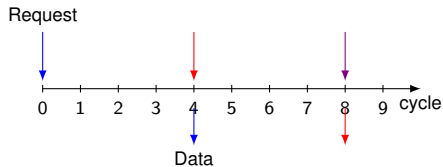


## Interleaved memory

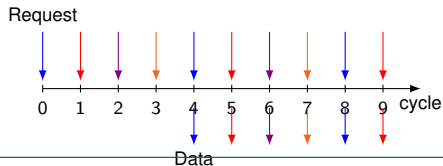
- Many physical memory units interleaved by the memory bus ;
- Number of physical memory units  $\equiv$  number of ways ;
- Number of contiguous bytes in a unique physical memory  $\equiv$  width of way ;
- Quadratic cost to build relative to number of memory units !

# Interleaved memory access

## Classic memory access



## four ways interleaved memory access



# Cache memory

---

## Consequences of grid architecture of RAMS

Bigger is a memory, bigger is her grid, slower is the read and write access.

## Cache memory

- Fast small memory unit where one store temporary data ;
- When multiple access to a same variable **in a short time**, speedup the read or write access ;
- Cache memory managed by the CPU (but cache memory for GPU can be managed by the programmer) ;
- **Consequence** : To optimize his program, the programmer must know the strategies used by the CPU.

# Cache memory

---

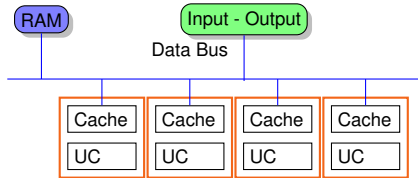
## CPU Strategy

- **Cache line** : store contiguous memory variables in cache (64 bytes on intel processor) ;
- **Associative memory cache** : Each cache memory address mapped on fixed RAM address (with a modulo).

## Consequences

- Better to have contiguous access in memory ;
- Better to use as soon as possible data stored in cache ;
- **Spatial and time localisation of data.**

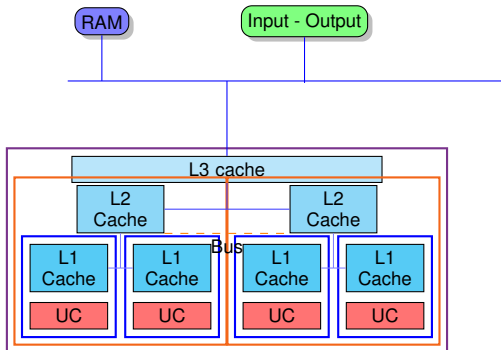
# memory organization on multi-processor computer



Data coherence between memory caches :

- An unique cache contains the datum : value is valid, none synchronization needed ;
- Datum shared with another memory caches : At each access, verify if the datum is modified by another core and write as invalid when modifying his value ;
- Modify the value in the cache : value now not valid in RAM : update the value in RAM if another core reads the value ;
- Value is invalid for cache. The next read of this value must access of the value in RAM.

# Many cores cache organization



Same issue with cache coherency, but need too coherence of data between cache levels. Complexity raises with the number of cache levels.

# Tools for shared memory computation

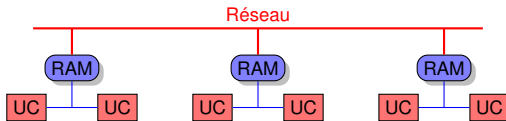
Many tools can be used to use many "threads" and the synchronization in memory. The most used are :

- OpenMP : Compilation directives and small C API (`#pragma`);
- Standard library in C++ (threads, asynchronous functions);
- TBB (Threading Building Block library, Intel) : Open source library from Intel

But memory conflict access must be cared for the programmer :

- When a thread writes a datum and some other threads read simultaneously the same datum ;
- When some thread writes at the same datum ;
- Doesn't rely on the instruction order in the program ! (out-of-order optimization by compiler and processor).

# Distributed memory



- Each computing unit can read/write on local ram : the set containing the computing unit and the ram is called **Computing node** ;
- The data are exchanged between computing nodes through a specialized bus or specific ethernet link ;
- On ethernet link, it's the responsibility of the programmer to exchange explicitly the data between computing nodes ;
- Need specifics efficiency algorithms and a library ;
- Possible to compute on many thousands of computing cores ;
- Only limited by electricity consumption (linear cost) ;



# Distributed parallelism context

All libraries managing the distributed parallel computation provide similar functionalities.

## Running an distributed parallel application

- An application is provided to the user to run his application(s) on a wanted number `nbp` of computing nodes (given when running the application) ;
- The computing nodes where the application(s) is launched is defined by default or in a file provided by the user ;
- The default output for all processes are the terminal output from which was launched the application ;
- A communicator (defining a set of processes) is defined by default including all launched processes (`MPI_COMM_WORLD`) ;
- The application gives an unique number for each process in a communicator (numbering from zero to `nbp-1`) ;
- All processes terminates the program at the same time ;

# Managing the context in your program

- Call initialization of parallel context before using other function in the library (MPI\_Init);
- Get the number of processes contained by the communicator (MPI\_Comm\_size);
- Read the rank of the process inside the communicator (MPI\_Comm\_rank);
- After calling the last library function, call the termination of parallel context to synchronize processes (MPI\_Finalize, if not done, crash your program).

```
#include <mpi.h>
int main(int nargs, char const* argv[])
{
    MPI_Comm commGlob;
    int nbp, rank;
    MPI_Init(&nargs, &argv); // Initialization of the parallel context
    MPI_Comm_dup(MPI_COMM_WORLD, &commGlob); // Copy global communicator in own communicator;
    MPI_Comm_size(commGlob, &nbp); // Get the number of processes launched by the used;
    MPI_Comm_rank(commGlob, &rank); // Get the rank of the process in the communicator commGlob.
    ...
    MPI_Finalize(); // Terminates the parallel context
}
```

# Point to point data exchange

A process send some data in a message to another process which receives this message.

## Constitution of a data message to send

- The communicator used to send the data ;
- The memory address of the contiguous data to send ;
- The number of data to send ;
- The type of the data (integer, real, user def type, and so.) ;
- The rank of destination process ;
- A tag number to identify the message

## Constitution of a data message to receive

- The communicator used to receive the data ;
- A memory address of a buffer where store the received data ;
- the number of data to receive ;
- The type of the data (integer, real, user def type, and so.)
- The rank of the sender process (can be any process) ;
- A tag number to identify the message (can be any tag if needed)
- Status of the message (receive status, error, sender, tag) ;

```
if (rank == 0) {  
    double vecteur[5] = { 1., 3., 5., 7., 22. };  
    MPI_Send(vecteurs, 5, MPI_DOUBLE, 1, 101, commGlob);  
} else if (rank==1) {  
    MPI_Status status;  
    double vecteurs[5];  
    MPI_Recv(vecteurs, 5, MPI_DOUBLE, 0, 101, commGlob, &status);  
}
```

# Interlocking

## Definition

- Interlocking is a situation where many processes are waiting each other for a infinite time to complete their messages ;
- By example, process one wait to receive a message from 0 and 0 wait to receive a message from 1 ;
- Or process zero sends a message to 1 and process one wait a message from 0 but with wrong tag !
- Sometimes, interlocking can be very hard to find !
- **Rule of thumb** : Be careful for each send to have a corresponding receive with right tag and expeditor.

```
if (rank==0)
{
    MPI_Recv(rcvbuf, count, MPI_DOUBLE, 1, 101, commGlob, &status);
    MPI_Send(sndbuf, count, MPI_DOUBLE, 1, 102, commGlob);
}
else if (rank==1)
{
    MPI_Recv(rcvbuf, count, MPI_DOUBLE, 0, 102, commGlob, &status);
    MPI_Send(sndbuf, count, MPI_DOUBLE, 0, 101, commGlob);
}
```

# Interlocking (more complicated cases)

```
MPI_Comm_rank(comm, &myRank ) ;
if (myRank == 0 )
{
    MPI_Ssend( sendbuf1, count, MPI_INT, 2, tag, comm);
    MPI_Recv( recvbuf1, count, MPI_INT, 2, tag, comm, &status);
}
else if ( myRank == 1 )
{
    MPI_Ssend( sendbuf2, count, MPI_INT, 2, tag, comm);
}
else if ( myRank == 2 )
{
    MPI_Recv( recvbuf1, count, MPI_INT, MPI_ANY_SOURCE, tag, comm,
              &status );
    MPI_Ssend( sendbuf2, count, MPI_INT, 0, tag, comm);
    MPI_Recv( recvbuf2, count, MPI_INT, MPI_ANY_SOURCE, tag, comm,
              &status );
}
```

# Blocking and non blocking message

---

## Definition

- Blocking message : Wait the complete reception of the message before returning from the function ;
- Non blocking message : Post the message to send or receive and return from the function immediatly !
- The status of non blocking message is update in a request struct. (not yet send/recv, sending/receiving or send/received)
- Allows to test or wait for the message to be completed.

## When use no blocking message ?

- When one can compute using other data during messages exchanges to hide message exchange cost ;
- To simplify algorithms to ensure none interlocking situations.

# Example using non blocking message

```
MPI_Request req;
if (rank == 0)
{
    double vecteur[5] = { 1., 3., 5., 7., 22. };
    MPI_Isend(vecteurs, 5, MPI_DOUBLE, 1, 101, commGlob, &req);
    // Some compute with other data can be executed here!
    MPI_Wait(req, MPI_STATUS_IGNORE);
}
else if (rank==1)
{
    MPI_Status status;    double vecteurs[5];
    MPI_Irecv(vecteurs, 5, MPI_DOUBLE, 0, 101, commGlob, &req);
    int flag = 0;
    do {
        // Do computation while message is not received on another data
        MPI_Test(&req, &flag, &status);
    } while(flag );
}
```

# A scheme to avoid interlocking situations

## The scheme for all processes

- First do receptions in non blocking mode ;
- Second, do send in blocking mode (or non blocking mode if you want to overlay message cost with computing)
- Third, synchronize yours receptions (waiting for completion or testing to overlay message cost with computing).

```
MPI_Request req; MPI_Status status;
if (rank==0)
{
    MPI_Irecv(rcvbuf, count, MPI_DOUBLE, 1, 101, commGlob, &req);
    MPI_Send(sndbuf, count, MPI_DOUBLE, 1, 102, commGlob);
    MPI_Wait(&req, &status);
}
else if (rank==1)
{
    MPI_Irecv(rcvbuf, count, MPI_DOUBLE, 0, 102, commGlob, &req);
    MPI_Send(sndbuf, count, MPI_DOUBLE, 0, 101, commGlob);
    MPI_Wait(&req, &status);
}
```



# Bufferized or non bufferized messages

## Bufferized messages

- A non blocking send is copied in a buffer before to be send ;
- ⇒ After calling a non blocking send, the user can modify the send data without changing the values to send ;
- It's the default behavior when we send a message
- But the copy of the data in a buffer has a memory and cpu cost !
- We can to call send functions which doesn't copy the data in a buffer ;
- It's the responsibility of the user to avoid to change data before the synchronization of the message !

# Example of non bufferized messages

```
std::vector<double> tab{3.,5.,7.,11.};  
// Blocking non bufferized send  
MPI_Ssend(tab.data(), tab.size(),  
          MPI_DOUBLE, 1, 104, commGlob);  
// OK, tab is sendd when we  
//   modify the buffer  
tab[3] = 13.;
```

Right example

```
std::vector<double> tab{3.,5.,7.,11.};  
// Non blocking non bufferized send  
MPI_Issend(tab.data(), tab.size(),  
           MPI_DOUBLE, 1, 104, commGlob,&request);  
MPI_Wait(&request, &status);  
// OK, tab is sendd when we  
//   modify the buffer  
tab[3] = 13.;
```

Right example

```
std::vector<double> tab{3.,5.,7.,11.};  
// Non blocking non bufferized send  
MPI_Issend(tab.data(), tab.size(),  
           MPI_DOUBLE, 1, 104, commGlob,&request);  
// ERROR, we modify the buffer before  
//   than the tab is sendd !  
tab[3] = 13.;
```

Wrong example

# Collective messages

---

## What is the collective communication

- Broadcast data from one process to all processes ;
- Scatter data from one process to all processes ;
- Gather data from all processes to one process ;
- Reduce data (with arithmetic operation) from all processes to one/all processes
- Scan data (with arithmetic operation) from all processes to all processes
- All to all broadcast/scatter data

## Why collective communication

- Point to point communication is enough for all algorithms !
- But, for some parallel operations (broadcasting, reduction, scattering), the optimal algorithm depends on net topology ;
- Distributed parallel libraries provides collective communication which are optimized for all net topology ;

# Distributed parallel rules

---

- Ethernet data exchange is very slow comparing to memory access : **limit as possible the data exchanges** ;
- To hide data exchange cost, it's better to compute some values during the exchange of data : **prefer to use non blocking message exchange**
- Each message has an initial cost : **prefer to regroup data in a temporary buffer if needed** ;
- All processes quit the program at the same time : **Try to balance the computing load of computing nodes** ;

# Disponible tools

---

- Program ethernet layers (for specialists only !);
- Use dedicated library (MPI, PVM, ...)

In all cases, data exchanges messages must call explicitly, calling functions provided by the library ;

Better to think his software in parallel at the beginning of the project.