

### Plan Formation Démarrage Rapide

- Introduction : Mésocentre CALMIP
  - **Concepts fondamentaux**
    - **Introduction à l'Architecture des systèmes HPC**
      - *Système à mémoire partagée / distribuée*
      - *Architecture Processeurs/ Accélérateurs*
      - *Présentation système de Calcul CALMIP : EOS*
      - *Visite salle Machine*

### Scientific computation

numerical simulation thanks to Scientific computation :

- Quantum Chemistry: Gaussian, deMon
- Atomic scale Material : VASP, Wien2K, CP2K
- Molecular Dynamic : Amber, GROMACS, ...
- Genomic : Blast, vina, (CHDB)
- Crash : LS-Dyna, Radioss
- Fluid Dynamic : Jädin, Neptune\_CFD, OpenFOAM, Fluent, STARCCM+
- Mechanics : ABAQUS, GETFEM++

Efficient (optimum?) execution of computer program  
 Time to solution !!!

Scientific Computation : transverse (all disciplines)....

Discretisation : spatial (FEM, Finite Volume, time (Explicit, implicit scheme, ...)  
 Numerical Methods (Solve Linear Systems, Algebra, ...)

Computer program ↔ High Performance Computing ↔ Mathematics

**Top500 November 2014** **Calcul Haute Performance : TOP 500 List**

RANK	SITE	SYSTEM	CORES	RMAX (TFLOP/S)	RPEAK (TFLOP/S)	POWER (KW)
1	National Super Computer Center in Guangzhou China	Tianhe-2 (MilkyWay-2) - TH-1VB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 3151P NUDT	3,120,000	33,862.7	54,902.4	17,808
2	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560,640	17,590.0	27,112.5	8,209
3	DOE/NNSA/LLNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1,572,864	17,173.2	20,132.7	7,890
4	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect Fujitsu	705,024	10,510.0	11,280.4	12,660
5	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786,432	8,586.6	10,066.3	3,945
6	Swiss National Supercomputing Centre (CSCS) Switzerland	Piz Daint - Cray XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect, NVIDIA K20x Cray Inc.	115,984	6,271.0	7,788.9	2,325
7	Texas Advanced Computing Center/Univ. of	Stampede - PowerEdge C8220, Xeon	442,462	5,168.1	8,520.1	4,510
26	Grand Equipement National de Calcul Intensif - Centre Informatique National de l'Enseignement Suprieur (GENCI-CINES) France	Occigen - bulx DLC, Xeon E5-2690v3 12C 2.6GHz, Infiniband FDR Bull SA	50,544	1,628.8	2,102.6	935
7	United States	1.600GHz, Custom Interconnect IBM				
10	Government United States	Cray CS-Storm, Intel Xeon E5-2660v2 10C 2.2GHz, Infiniband FDR, Nvidia K60 Cray Inc.	72,800	3,577.0	6,131.8	1,499

**Moore's Law ?**

<http://www.calmip.univ-toulouse.fr/>

**Top500 November 2015**

RANK	SITE	SYSTEM	CORES	RMAX (TFLOP/S)	RPEAK (TFLOP/S)	POWER (KW)
1	National Super Computer Center in Guangzhou China	Tianhe-2 (MilkyWay-2) - TH-1VB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 3151P NUDT	3,120,000	33,862.7	54,902.4	17,808
2	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560,640	17,590.0	27,112.5	8,209
3	DOE/NNSA/LLNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1,572,864	17,173.2	20,132.7	7,890
4	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect Fujitsu	705,024	10,510.0	11,280.4	12,660
5	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786,432	8,586.6	10,066.3	3,945
6	DOE/NNSA/LANL/SLNL United States	Trinity - Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Aries interconnect Cray Inc.	301,056	8,100.9	11,078.9	
7	Swiss National Supercomputing Centre (CSCS) Switzerland	Piz Daint - Cray XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect, NVIDIA K20x Cray Inc.	115,984	6,271.0	7,788.9	2,325
8	HLRS - Höchstleistungsrechenzentrum Stuttgart Germany	Hazel Hen - Cray XC40, Xeon E5-2680v3 12C 2.5GHz, Aries interconnect Cray Inc.	185,088	5,640.2	7,403.5	
9	King Abdullah University of Science and Technology Saudi Arabia	Shaheen II - Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Aries interconnect Cray Inc.	196,608	5,537.0	7,235.2	2,834
10	Texas Advanced Computing Center/Univ. of Texas United States	Stampede - PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Infiniband FDR, Intel Xeon Phi SE10P Dell	442,462	5,168.1	8,520.1	4,510

<http://www.calmip.univ-toulouse.fr/>

### Calcul Intensif et Système HPC

➤ **Calcul Intensif : Principes des systèmes HPC**

- ❑ Hardware et Software => Performance Calcul Flottants
  - ❑ Flop/s + Mémoire (RAM et espace fichier)
  - ❑ Stockage, I/O
  - ❑ €++ => Mutualisation
- ❑ Plusieurs Utilisateurs d'un même serveur
  - ❑ Partage des ressources cohérent : Règles Utilisation
  - ❑ OS performant : Multi-applications, Multi-User
- ❑ Serveur Totalement Dédié au Calcul
  - ❑ Applications Scientifiques Calcul uniquement
  - ❑ Sauvegarde
  - ❑ Espace Fichier / Stockage
  - ❑ Accès distant
- ❑ Flop : floating operation (mult, add) => opération sur les nombres à virgule flottantes (nombre réels)
  - ❑ 3,14159
  - ❑ -6,8675456342 E+08
- ❑ contrainte d'hébergement lourdes :
  - ❑ Electricité, (secouru)
  - ❑ Refroidissement
  - ❑ Poids
  - ❑ Sécurité
  - ❑ ...

Espace Clément Ader  
<http://www.calmip.univ-toulouse.fr/>

### Panorama Systèmes HPC

**Machine à Mémoire Partagée :**

- Multiprocesseurs
- Un seul espace d'adressage
- Mémoire partagée

**Symetric Multi-Processing SMP**

**UMA**  
Accès Mémoire Uniforme

**NUMA**  
Accès Mémoire Non Uniforme

---

**PROGRAMME**

↑

↓

**Machine à Mémoire Distribuées :**

- Multi-Ordinateurs
- Espace d'Adressage Multiple

**Massively Parallel Processing MPP Clusters**

**NORMA**  
no-remote memory access

Espace Clément Ader  
<http://www.calmip.univ-toulouse.fr/>

### Système d'Exploitation : Définition

**Système d'Exploitation (ou Operating System « OS » en anglais) :**  
 Un ensemble de logiciels ou programmes qui permet d'unifier les ressources matérielles, pour qu'elles soient utilisables par l'utilisateur.  
 Exemple : Windows, Linux, MacOSX, AIX (IBM), HPUX (HP), SOLARIS (SUN)

The diagram illustrates the layers of an operating system. At the bottom is 'Matériel informatique' (computer hardware). Above it is the 'Système d'exploitation' (operating system). The top layer is 'Application', which is used by the 'Utilisateur' (user). Bidirectional arrows indicate the flow of data and control between these layers.

(Schéma Operating system Wikipedia)

---

### UMA Architecture (Shared Memory)

- Machine side: SMP Symmetric MultiProcessor (SMP)
  - Bus Interconnexion between memory and processors
  - Central memory and I/O : shared by all processors
  - Processors access to the same memory (address space)
- User side:
  - A single machine (single OS) – several processors – one single space memory address
  - How to program : extension of sequential programming

The diagram shows a shared memory architecture. A large yellow box labeled 'memory' is connected to a 'BUS interconnect' (yellow bar). Multiple 'Processor' units are connected to this bus. Each processor contains a 'Cache', a 'Cache Coherency Unit', a 'Register File', and 'Functional Unit (mult, add)'. The 'Cache' and 'Cache Coherency Unit' are connected to the 'BUS interconnect'. The 'Register File' and 'Functional Unit' are connected to the 'Cache'. The 'OS' label is positioned at the end of the bus.

---

### UMA Architecture (Shared Memory) - Multithreading

> Parallel Programming with Shared Memory Architecture : OpenMP

**AUTOMATIC LOOP PARALLELISATION !!!!**

```

!SOMP DO PARALLEL
do i = 1, n
a(i) = 92290. + real(i) ;
end do
!SOMP END DO
    
```

Automatic : spread loop's iterations on cores

11/02/16 Page 9 <http://www.calmip.univ-toulouse.fr/>

### UMA Architecture

**Memory access:**

- Concurrent access to central memory => bottleneck
  - Time access increase
- Increase size (and so level) of caches

**Consequence : few number of processor**

•Another paradigma/option : distribute memory ?

11/02/16 Page 9 <http://www.calmip.univ-toulouse.fr/>

### Distributed memory(NORMA)

- Processor and memory tightly interconnected
  - MPP : Massively Parallel Processing
  - Cluster : machines(comput nodes) interconnection

Architecture System Share

Architecture	Share
Cluster	85.2%
MPP	14.8%

Espace Clément Ader

<http://www.calmip.univ-toulouse.fr/>

### Distributed memory : multi-computer Architecture (Clusters)

- Machine side:
  - Massive technology
  - Process access to its own (local) memory space
  - Interconnect nodes :
    - Like internet (ethernet)...
    - need much faster (bandwidth and latency)
    - process to process communication

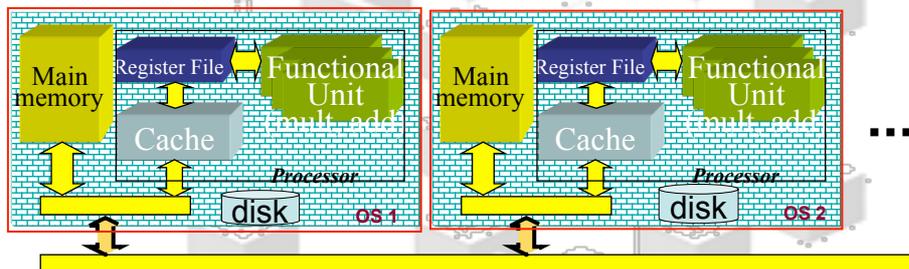
Espace Clément Ader

<http://www.calmip.univ-toulouse.fr/>

### Distributed memory : multi-computer Architecture (Clusters)

- User side:

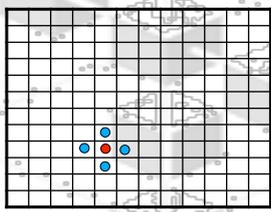
- $n$  different nodes ( $n$  OS) interconnected, 1 (or +) processor per node.
- Parallel programming  $\rightarrow$  Message Passing Interface (exchange messages, work done by developer ...you?)
- Need efficient tools to properly access computing resources



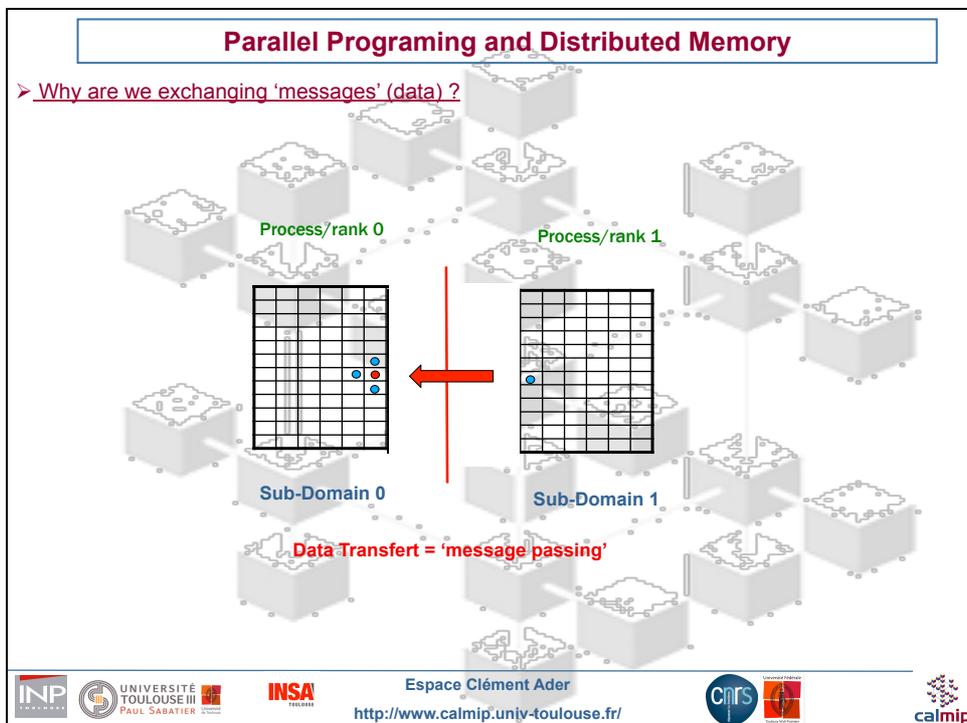
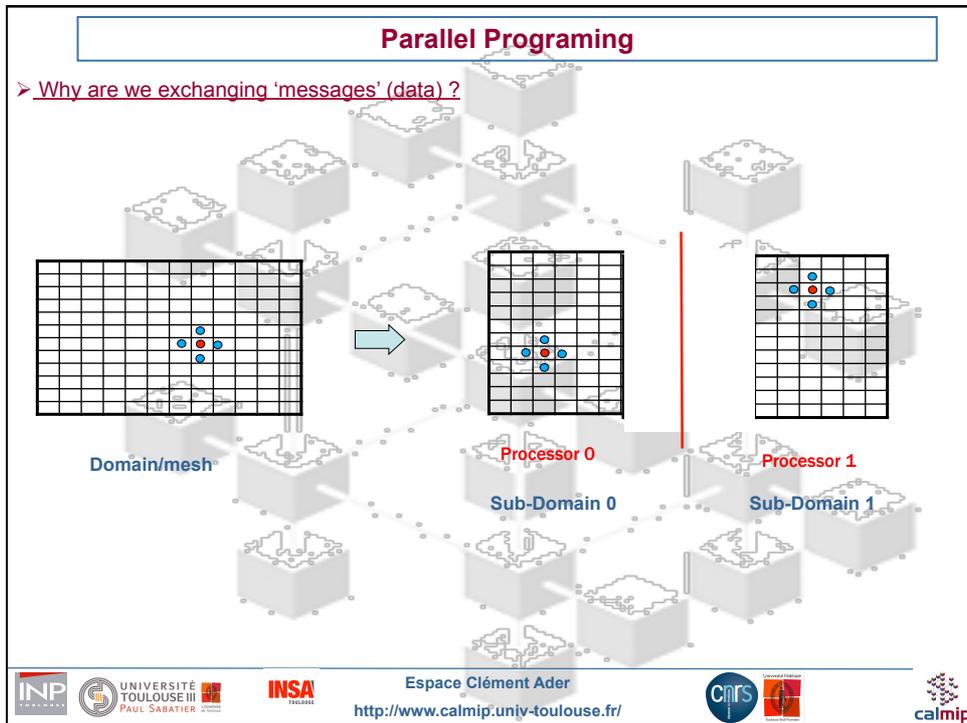
### Parallel Programming

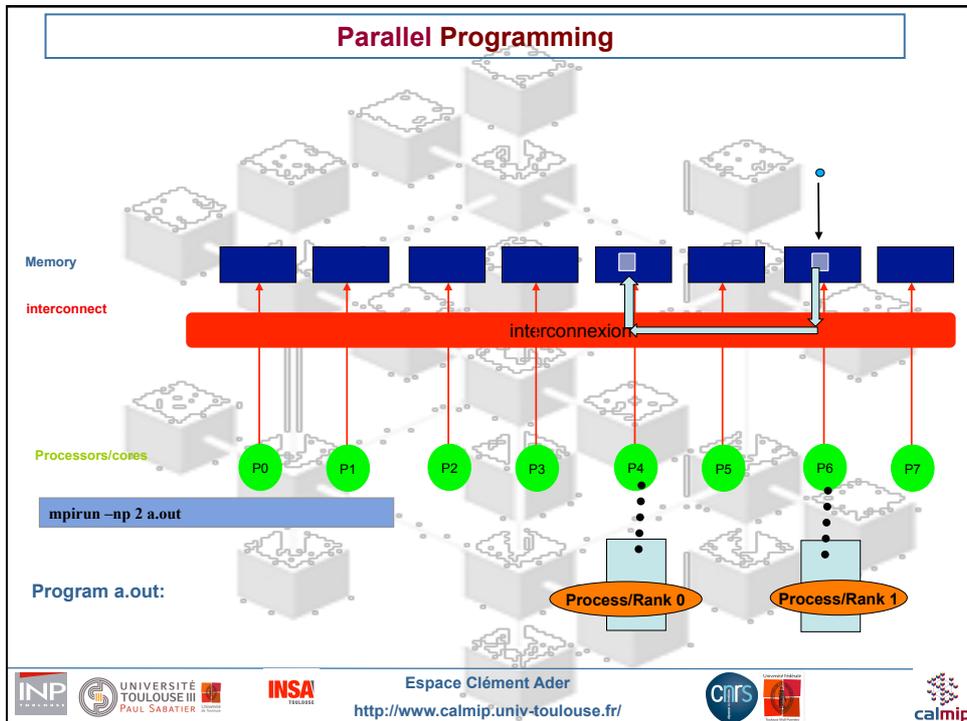
> Why are we exchanging 'messages' (data) ?

Compute the Red (i,j) point thanks to values of blue points (neighbors)



Domain/mesh/grid





### Parallel Programming

➤ SEND/RECV routine

```

program point_to_point
USE MPI

call MPI_INIT(code)
call MPI_COMM_RANK(MPI_COMM_WORLD,rank)

if (rank == 6) then
    value=1000
    call MPI_SEND(value,1,MPI_INTEGER,4,...,MPI_COMM_WORLD)
elseif (rank == 4) then
    call MPI_RECV(value,1,MPI_INTEGER,6,...,MPI_COMM_WORLD)
    print *, 'I, process', 4, 'I recieved', value, ' from process 6'
end if
call MPI_FINALIZE(code)
end program point_to_point
    
```

Espace Clément Ader  
<http://www.calmip.univ-toulouse.fr/>

### Architecture à mémoire partagée physiquement distribuée

**Machine NUMA : Non-Uniform Memory Access**

- Mixage SMP/cluster
- Souplesse d'utilisation et performance
- Dépendance « spatiale » par rapport aux ressources utilisées

**Réseau d'interconnexion**

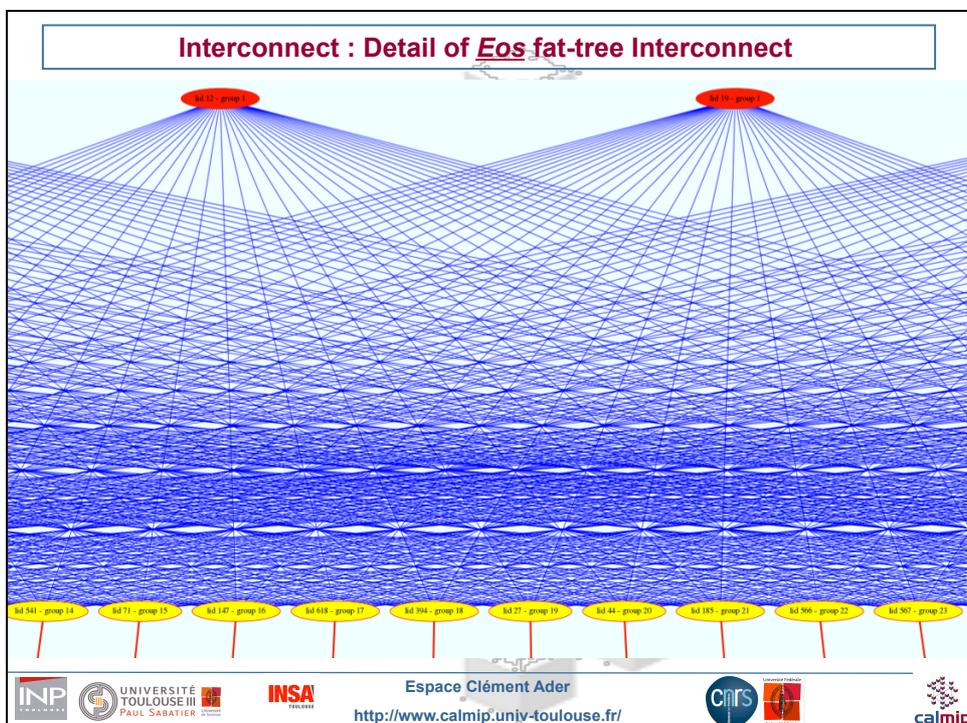
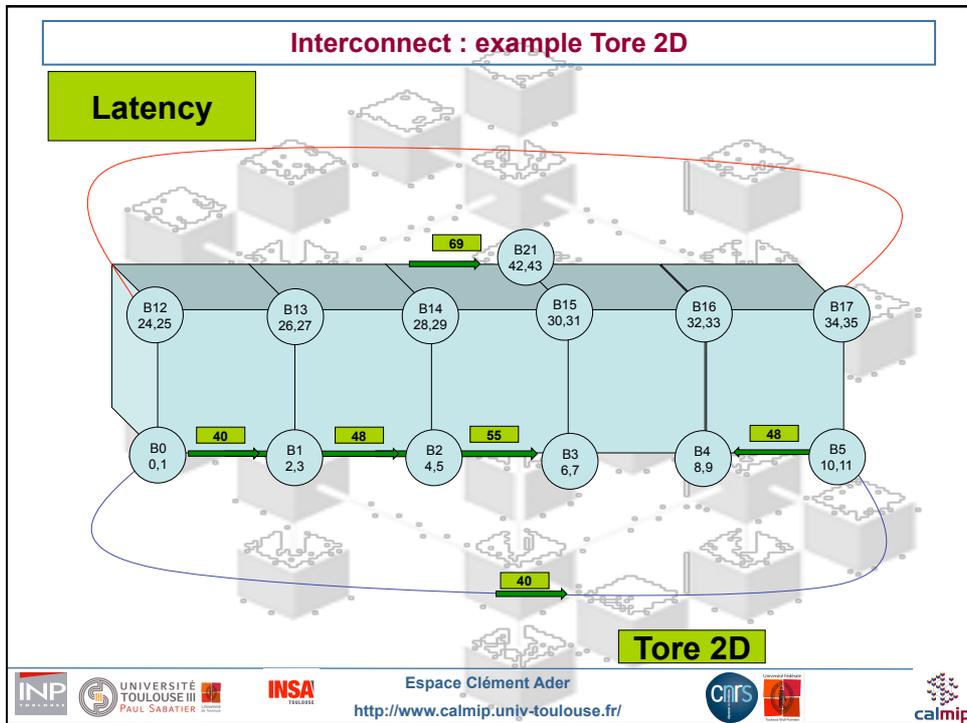
<http://www.calmip.univ-toulouse.fr/>

### Interconnexion : key of parallel machine

- ❑ Interconnexion :
  - ❑ Latency : how much time to be connected ? Order of microsecond
  - ❑ bandwidth (throughput): rate of data transfer ? Mbytes/sec
  - ❑ Topology : how many path from a point to another
- ❑ In parallel machine hardware (processor / memory) have to be connected
- ❑ specific (fast) protocols
  - ❑ (cluster) infiniband, proprietary
  - ❑ topology:
    - ❑ ring, hypercube, Torus, fat-tree, ...

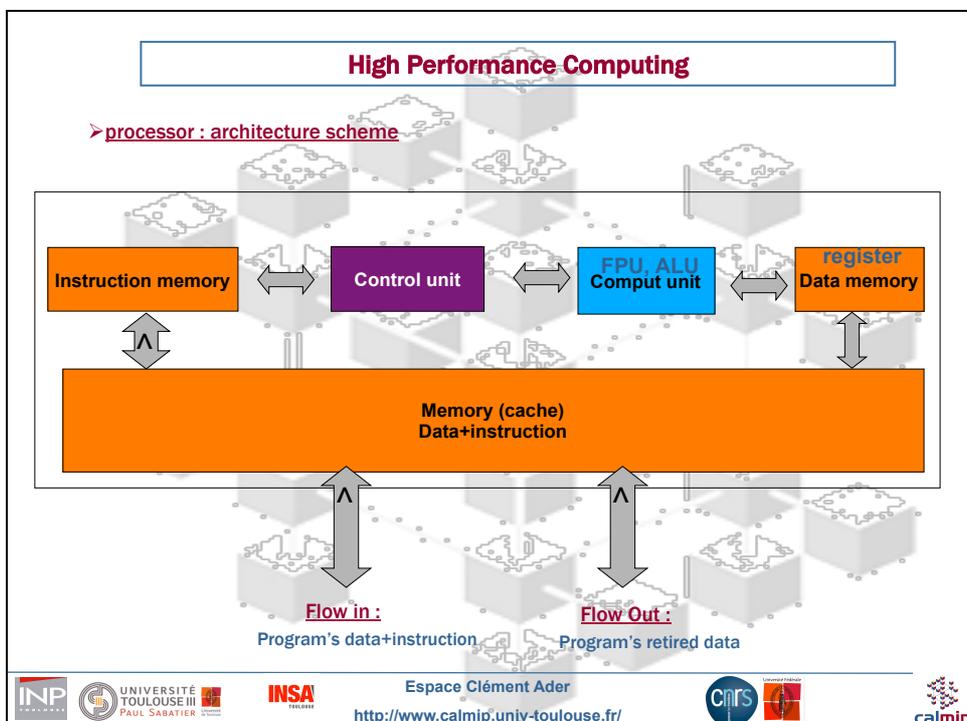
1 Mo/s	1 Megaoctet/s	10 <sup>6</sup> octets/sec
1 Go/s	1 Gigaoctet/s	10 <sup>9</sup> octets/sec
1 To/s	1 Téraoctet/s	10 <sup>12</sup> octets/s

<http://www.calmip.univ-toulouse.fr/>





- Introduction : Mésocentre CALMIP
  - **Concepts fondamentaux**
    - **Introduction à l'Architecture des systèmes HPC**
      - *Système à mémoire partagée / distribuée*
      - **Architecture Processeurs/ Accélérateurs**
      - *Présentation système de Calcul CALMIP : EOS*
      - *Visite salle Machine*



### High Performance Computing

➤ processor : architecture scheme

➤ processor cycle

- Clock frequency = number of pulse per second
- 1,5 Ghz ⇒ 1,5 Billion cycles per second
- 1 cycle : perform one « atomic » processor instruction

1,5 Ghz ⇒ 1 cycle = 0,6 ns

---

### High Performance Computing

#### Moore's Law

Year	Processor	Transistors
1971	4004	~23,000
1974	8080	~60,000
1978	8086	~290,000
1982	80286	~1,200,000
1985	80386	~2,750,000
1989	Pentium Processor	~3,100,000
2000	Micro 2000	~90,000,000

- ❑ « something » double each 18 month
- ❑ « something » = transistors
- ❑ « something » related to « performance »
- ❑ this is an exponential Law
  - ❑ could it last ? (True for 30 years)
  - ❑ empirical anyway

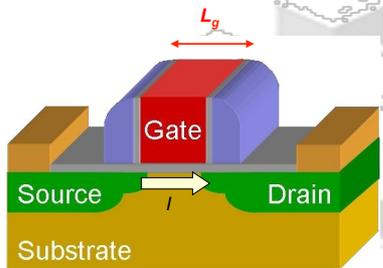
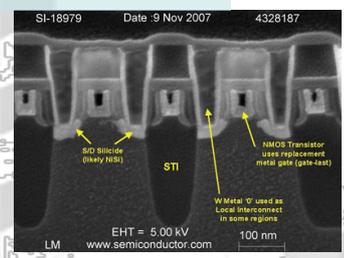
**Moore(Intel), Moore'Law is not a Physical law, but depend a lot on ... Physics!**

<http://www.nature.com/nature/journal/v512/n7513/full/nature13570.html>

---

### High Performance Computing

▣ A matter of Engraving a piece of Silicon :

▣ Energy Problem

$$W = CV^2 \quad \text{Power} = W * F$$

Consequence : End of the Race for Frequency





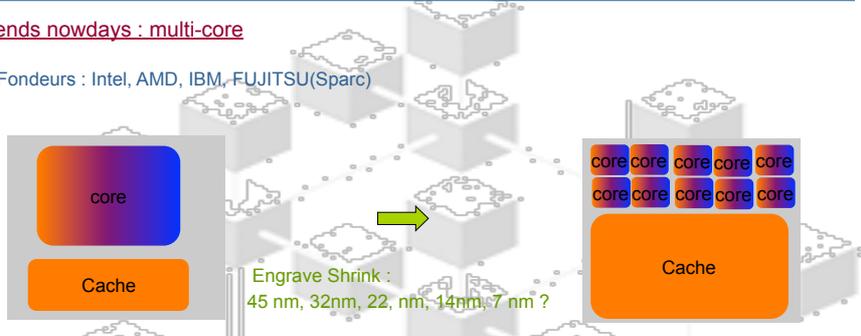
Espace Clément Ader  
<http://www.calmip.univ-toulouse.fr/>




### Processors Architecture HPC

➤ Trends nowadays : multi-core

•Fondeurs : Intel, AMD, IBM, FUJITSU(Sparc)



Engrave Shrink : 45 nm, 32nm, 22, nm, 14nm, 7 nm ?

- 1 processeur mono-core
- 1 processeur ( or socket ) multi-core (multi = 2, 4, 6, 8, 10, 12, 20 )
- More 'RAW' power (x2, 4, 6, 8) => parallelism
- better ratio flop/watt, flop/m²
- same frequency or lower !





Espace Clément Ader  
<http://www.calmip.univ-toulouse.fr/>




### Processor Architecture : SIMD

Vector Processor

- From computing scalar to computing Vector
- SIMD : SINGLE INSTRUCTION MULTIPLE DATA
- each processor cycle  $n$  scalar retired operations

```
for (i = 0; i <= MAX; i++)
  c[i] = a[i] + b[i];
```

<b>a[i]</b>	a										
		a[+7]	a[+6]	a[+5]	a[+4]	a[+3]	a[+2]	a[+1]	a[+0]		
		b[+7]	b[+6]	b[+5]	b[+4]	b[+3]	b[+2]	b[+1]	b[+0]		
		c[+7]	c[+6]	c[+5]	c[+4]	c[+3]	c[+2]	c[+1]	c[+0]		

42 / 75    Software & Services Group, Energy Engineering Team    intel

Espace Clément Ader

<http://www.calmip.univ-toulouse.fr/>

### Processors Architecture : SIMD

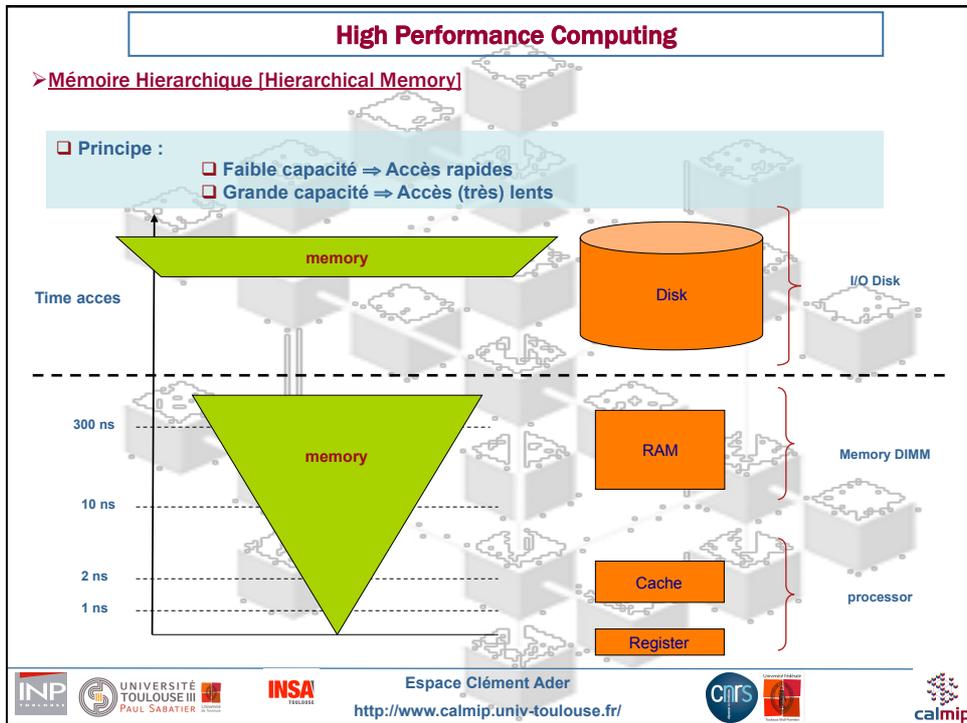
Processor Architecture

```

graph LR
    IM[Instruction memory] <--> CU[Control unit]
    CU <--> CU2[Comput unit FPU, ALU]
    CU2 <--> DM1[Data memory]
    CU2 <--> DM2[Data memory]
    CU2 <--> R[register]
    CU2 <--> DM3[Data memory]
    CU2 <--> DM4[Data memory]
    IM <--> MC[Memory cache Data+instruction]
    CU <--> MC
    CU2 <--> MC
    
```

Espace Clément Ader

<http://www.calmip.univ-toulouse.fr/>



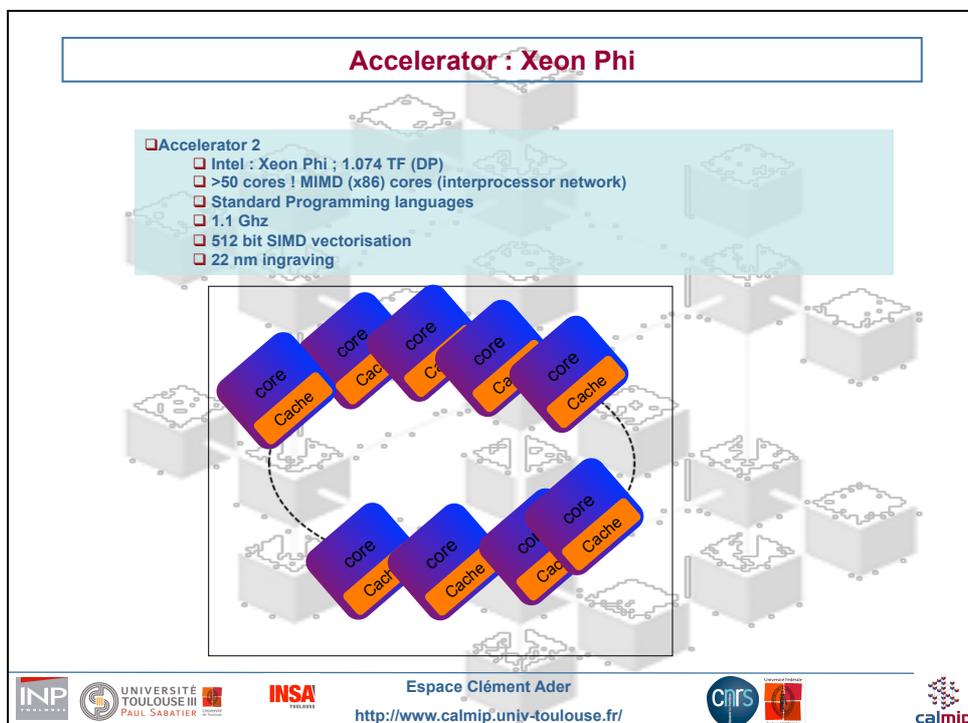
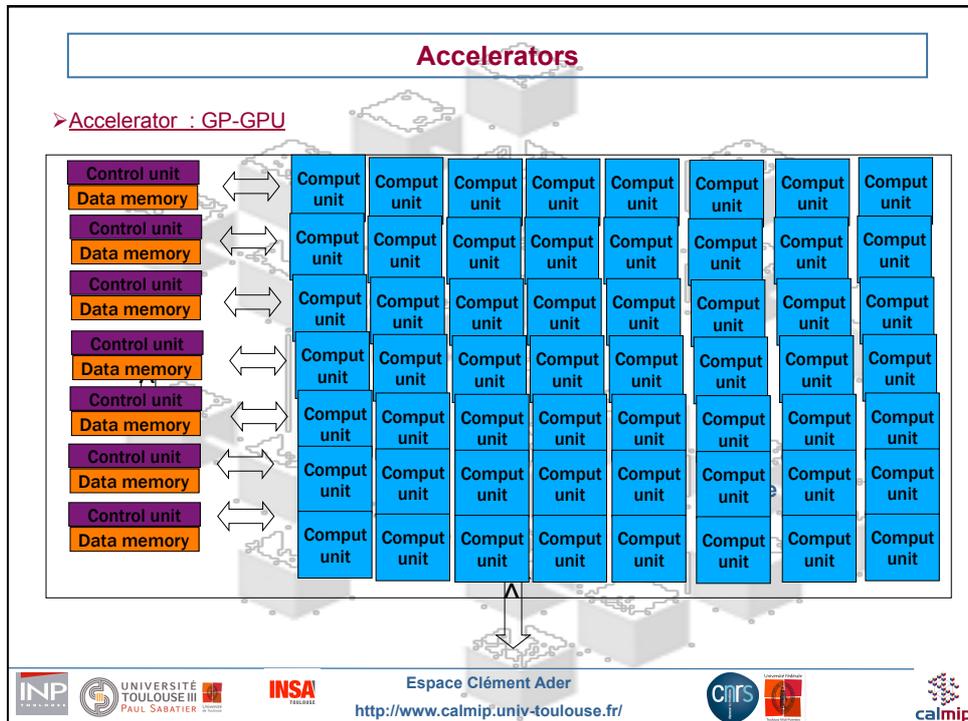
### Top500 Juin 2015 Calcul Haute Performance : TOP 500 List

	RANK	SITE	SYSTEM	CORES	RMAX (TFLOP/S)	RPEAK (TFLOP/S)	POWER (KW)
#1 2015	1	National Super Computer Center in Guangzhou China	Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 3151P NUDT	3,120,000	33,862.7	54,902.4	17,808
#1 2014	2	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7, Opteron 6274.16C.2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560,640	17,590.0	27,112.5	8,209
#1 Nov. 2012	3	DOE/NNSA/LLNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1,572,864	17,173.2	20,132.7	7,890
#1 2012	4	RIKEN Advanced Institute for Computational Science (AICS) Japan	TK1 - Intel Xeon Phi 3151P, Tofu interconnect Fujitsu	705,024	10,510.0	11,280.4	12,660
#1 2011	5	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786,432	8,586.6	10,066.3	3,945
	6	Swiss National Supercomputing Centre (CSCS) Switzerland	Piz Daint - Cray XC30, Xeon E5-2670.8C 2.400GHz, Aries interconnect, NVIDIA K20x Cray Inc.	115,984	6,271.0	7,788.9	2,325
	7	King Abdullah University of Science and Technology Saudi Arabia	Shaheen II - Cray XC40, Xeon E5-2698v3 16C 2.30GHz, Aries interconnect Cray Inc.	196,608	5,537.0	7,235.2	2,834
	8	Texas Advanced Computing Center/Univ. of Texas United States	Stampede - PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Infiniband FDR, Intel Xeon Phi SE10P Dell	462,462	5,168.1	8,520.1	4,510

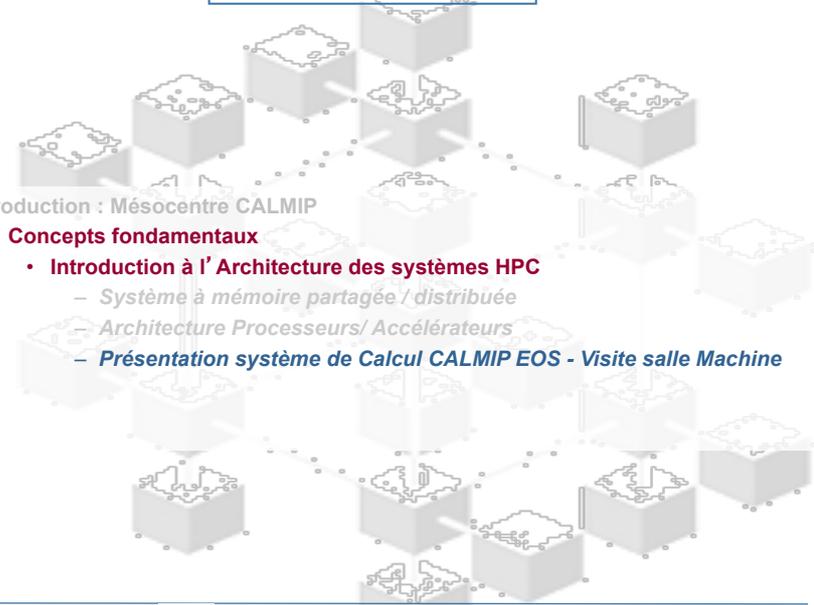
**Accélérateurs**

Nuclear Plant between 40 MW and 1450 MW

Espace Clément Ader  
<http://www.calmip.univ-toulouse.fr/>



**Plan Formation Démarrage Rapide:**



- Introduction : Mésocentre CALMIP
  - **Concepts fondamentaux**
    - **Introduction à l' Architecture des systèmes HPC**
      - *Système à mémoire partagée / distribuée*
      - *Architecture Processeurs/ Accélérateurs*
      - *Présentation système de Calcul CALMIP EOS - Visite salle Machine*

**(Brand new) CALMIP's Supercomputer : EOS**



- Bullx DLC Cluster
- Number of cores : 12240 [612 nodes]
- 39.1 TB distributed memory
- June 2014 : Reach #183@TOP500 ranking
- 274 TF Peak - Linpack (Rmax) : 255 TF

Page web associée : <http://www.calmip.univ-toulouse.fr/spip/spip.php?article388>




 Espace Clément Ader  
<http://www.calmip.univ-toulouse.fr/>



