




Multiprocessors

- One system --- many processors
 - parallelism

Jesús Labarta, MP, 2008



Paralelismo

■ Expectativa

- más currantes
 - ✓ menos tiempo
 - ✓ más trabajo en el mismo tiempo

..... lineal

Jesús Labarta, MP, 2008



Parallelism

■ Reality

- Work partitioning: sometimes
 - ✓ Lack of work: **dependencies** (sequential is an extreme case)
 - Communication / synchronization
 - ✓ Hard to equipartition: **Load imbalance**
 - ✓ **replication**
- Overhead
 - ✓ Work generation
 - ✓ Communication: volume/surface ratio
 - ✓ synchronization

Jesús Labarta, MP, 2008



Parallelism

■ Result

- Team delivers more than one individual
- **Not always linear**
- Sometimes (few) more than the sum of the parts

Jesús Labarta, MP, 2008



Multiprocesadores (2004)

■ Desde los dinosaurios



EARTH SIMULATOR
Earth Simulator Center
Kanazawa
NEC
Rmax: 35.86 TFlops

<http://www.top500.org/>



ASCI WHITE
LLNL
Livermore
IBM SP Power3
Rmax: 7.22 TFlops



TERASCALE SYSTEM
Pittsburgh
Supercomputer Ctr.
HP AlphaServer
Rmax: 4.46 TFlops



TERA
CEA
Bruyeres-le-Chatel
HP AlphaServer SC
Rmax: 3.98 TFlops



SP POWER3
NERSC/LBNL
Berkeley
IBM SP Power3
Rmax: 3.05 TFlops

■ Hasta los ratones (super)



Jesús Labarta, MP, 2008



Multiprocesadores (2005)

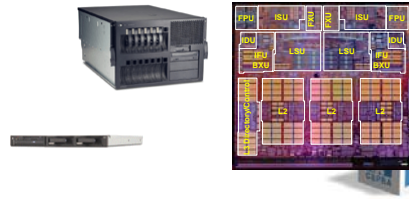
■ Desde los dinosaurios

- ... todos tuvieron sus 5 minutos de gloria



■ Hasta los ratones (super)

- ...



Jesús Labarta, MP, 2008

Multiprocessors (2008)

■ What changed?

- From elite: still there
- To “commodity”

■ Why?

- Necessity: Power

Jesús Labarta, MP, 2008



Multiprocessors

TOP 500

- <http://www.top500.org/>



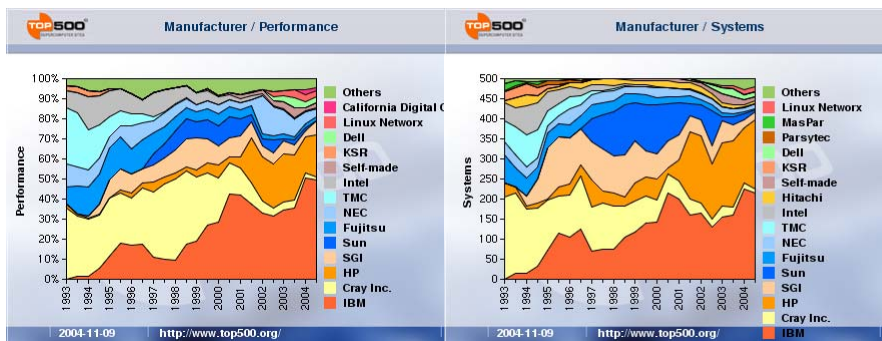
Jesús Labarta, MP, 2008



Multiprocessors

TOP 500

- <http://www.top500.org/> by manufacturer



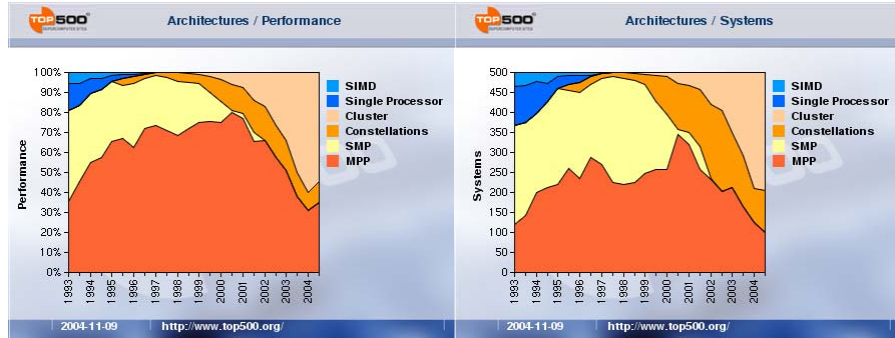
Jesús Labarta, MP, 2008



Multiprocessors

TOP 500

- <http://www.top500.org/> by architecture



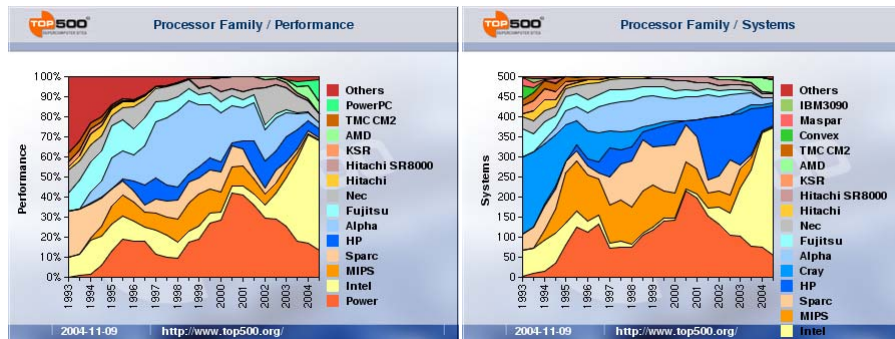
Jesús Labarta, MP, 2008



Multiprocessors

TOP 500

- <http://www.top500.org/> by processor



Jesús Labarta, MP, 2008



Scaling models

■ Scaling:

- Evolution of performance indices when increasing
 - ✓ Problem size
 - ✓ Processors
- Correlation on how to increase them?

Jesús Labarta, MP, 2008



Scaling models

■ Models:

- Problem Constrained scaling
 - ✓ Fixed problem size
- Time-constrained scaling
 - ✓ Keep constant total execution time
- Memory-constrained scaling
 - ✓ Keep constant memory used per processor

$$Speedup_{PC}(p) = \frac{Time(1)}{Time(p)}$$

$$Speedup_{TC}(p) = \frac{Work(p)}{Work(1)}$$

$$Speedup_{MC}(p) = \frac{IncreasedWork}{IncreasedTime}$$

Jesús Labarta, MP, 2008



Parallelism: useful

■ Capability computing

- Big problems
 - ✓ Not just time (CPUs)
 - ✓ Also Memory, I/O,...

■ Capacity computing:

- Throughput. Many problems

Jesús Labarta, MP, 2008



Multiprocesadores

■ Que tenemos nosotros (2004)

- | | |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none">● SGI Origin2000:<ul style="list-style-type: none">✓ 64 MIPS R10000● IBM SP:<ul style="list-style-type: none">✓ : 128 Power3✓ : 32 Power4● HP (compaq):<ul style="list-style-type: none">✓ : 12 Alpha 21264✓ : 16 Alpha 21264 | <ul style="list-style-type: none">● Cluster<ul style="list-style-type: none">✓ 64 Pentiums III✓ 16 Pentiums II● SMPs pequeño<ul style="list-style-type: none">✓ 4 Pentiums✓ 4 Itaniums |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

.... no estamos en el mapa ...pero se puede hacer mucho

Jesús Labarta, MP, 2008



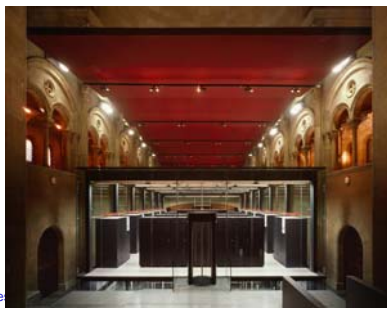
Multiprocesadores

■ Que tenemos nosotros (2005)

- Mare Nostrum:
 - ✓ ~5200 PPC970

.... cuartos en la foto (Nov 2004)

... octavos (Nov 2005)



...



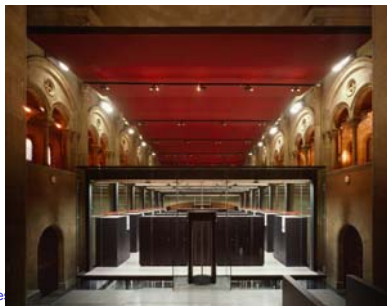
Multiprocessors

■ Que tenemos nosotros (2008)

- Mare Nostrum (2006)
 - ✓ ~10400 cores PPC970MX

.... 13th in the photo (Nov 2007)

- Altix:
 - ✓ ~128 cores Itanium



...



Multiprocessors

■ Our resources (2010) ??

- Mare Incognito:
 - ✓ ~ 600000 cores ??
 - ✓ Cell ??
 - ✓ ????

.... ????? in the photo (2011)

...

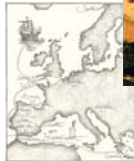
Jesús Labarta, MP, 2008



Mare Incognito

“Homogeneous” supercomputer based on the Cell processor

We believe it is possible to build
“cheap”, “efficient”, “not specific” and
“widely applicable”



We know it is risky

We have “a vision” of some relevant technologies to develop/
push forward

Many of them are not Cell specific and will be evaluated for other
architectures

An Opportunity

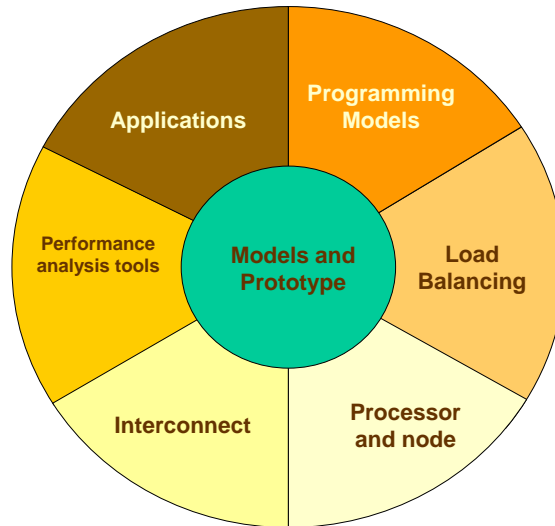
Integrate many BSC research lines and increase cooperation IBM

Influence the design and use of supercomputers in the future.

Jesús Labarta, MP, 2008



Project structure



Jesús Labarta, MP, 2008



Performance models

■ Model

- Relationship between
 - ✓ System parameters
 - ✓ Performance metrics
- Useful
 - ✓ Predict
 - Expectations, Reference
 - ✓ Understand behavior

Not necessarily bad when a model does not match reality...

... indication of unaccounted factor, not understood phenomena.

Jesús Labarta, MP, 2008



Performance models

■ Performance metrics

- Time $T(p)$
 - ✓ Speed-up $S(p) = T(1)/T(p)$
 - ✓ Efficiency $\eta(p) = S(p)/p$
- Throughput $??/T$ $?? \in \{\text{jobs, iterations, problems...}\}$

■ Other metrics

- Price/performance
 - ✓ MFLOP/\$
 - ✓ MFLOP/Watt
 - ✓ MFLOP/m²

Jesús Labarta, MP, 2008



Performance models

■ Amdahl's law

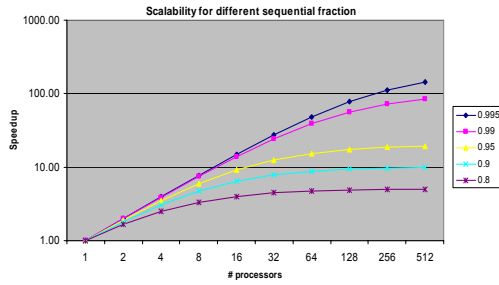
- $T(1) = T_{\text{seq}} + T_{\text{par}}$
- $T(P) = T_{\text{seq}} + T_{\text{par}}/P$
- $S(P) = 1/((1-f)+f/P)$
 - ✓ $f = T_{\text{par}}/T(1)$ parallel fraction
- Lo que no puede ser no puede ser y además es imposible

.... What is not parallelized ends up being the bottleneck

Jesús Labarta, MP, 2008



Speed-up: Amdahl's Law

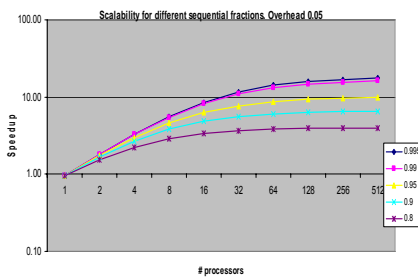


$$S(p) = 1 / (1 - f + (f / p))$$

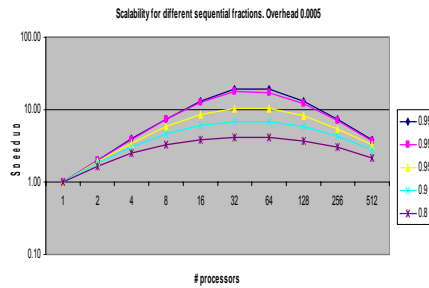
Jesús Labarta, MP, 2008



Speed-up: Overheads



$$S(p) = 1 / (1 - f + (f / p) + o)$$



$$S(p) = 1 / (1 - f + (f / p) + o * p)$$

Jesús Labarta, MP, 2008

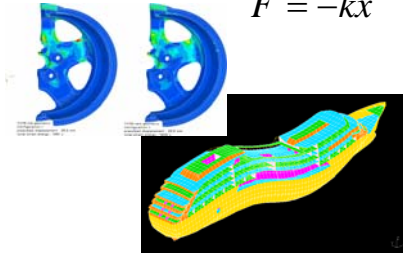


Laws of ...

■ Science

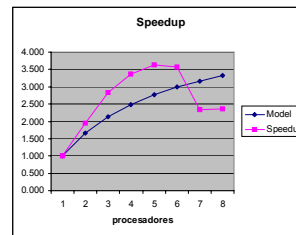
$$F = ma$$

$$F = -kx$$



■ Computer Science

$$S = \frac{1}{1 - f + \frac{f}{P}}$$



Jesús Labarta, MP, 2008



Architectures

■ Shared Memory

- SMP
- NUMA

■ Distributed Memory

■ Hierarchical mix

- SMP nodes @ Distributed memory

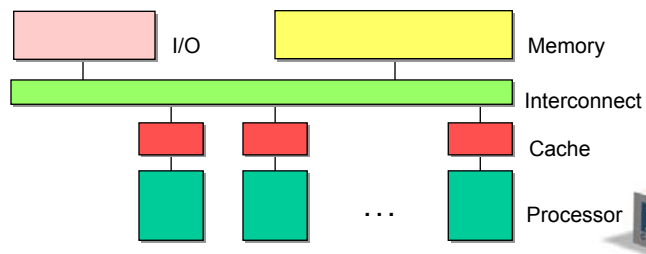
Jesús Labarta, MP, 2008



Shared Memory

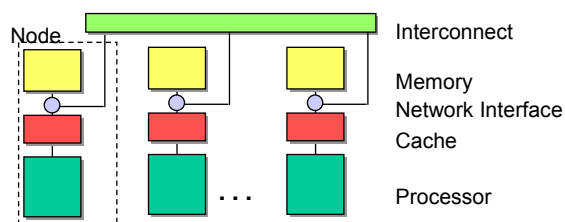
■ SMP (symmetric multiprocessor):

- UMA
- Issues
 - ✓ Interconnection network:
 - Aggregated bandwidth: bus, crossbar.
 - Latency
 - ✓ Cache: coherency, consistency



Jesús Labarta, MP, 2008

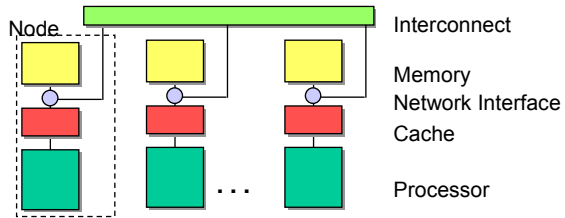
Distributed Shared Memory



- Load/stores: local/remote depending on physical address
- NUMA (non-uniform memory access): remote accesses may take 5-10 times longer

Jesús Labarta, MP, 2008

Distributed Shared Memory



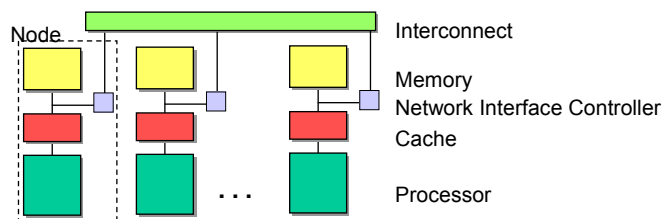
■ Issues:

- Interconnection network: latency, topology, bisection bandwidth
- Cache: coherency, consistency !!
- Locality of access

Jesús Labarta, MP, 2008



Distributed Memory

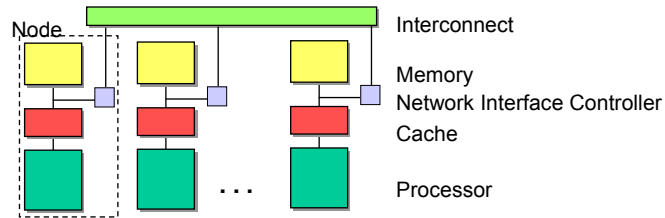


- Load/stores: local
- I/O: send/receive (local address + node addressing mechanism)
- Private / commodity networks

Jesús Labarta, MP, 2008



Distributed Memory



■ Issues

- Interconnection network:
 - ✓ latency, topology, bisection bandwidth
 - ✓ Injection mechanism

■ SMP nodes

Jesús Labarta, MP, 2008



Hierarchical

- Distributed memory
- Each node shared memory
 - ~ SMP several chips
- Each chip multicore/manycore

- Homogeneous/heterogeneous

Jesús Labarta, MP, 2008



Modelos de programación

■ Mecanismos disponibles al programador para expresar la estructura lógica de un programa

■ Influye

- Complejidad del programa
 - ✓ Costo de desarrollo
 - ✓ Legibilidad. Costo de mantenimiento
- Rendimiento
 - ✓ Influenciado por el modelo
 - ✓ por la implementación del modelo
 - ✓ Por la estructura de paralelización

Jesús Labarta, MP, 2008



Back to Babel?

Book of Genesis

"Now the whole earth had one language and the same words" ...

..."Come, let us make bricks, and burn them thoroughly."...

..."Come, let us build ourselves a city, and a tower with its top in the heavens, and let us make a name for ourselves"...

And the LORD said, "Look, they are one people, and they have all one language; and this is only the beginning of what they will do; nothing that they propose to do will now be impossible for them. Come, let us go down, and confuse their language there, so that they will not understand one another's speech."

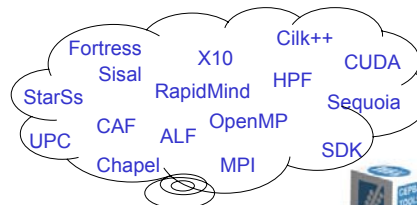


The computer age

Fortran & MPI



++



Jesús Labarta, MP, 2008



Back to Babel?

■ Result?

- Dispersion
- Pain
- Thousands of years from the first confusion of tongues to the second

■ Current situation

- We are at a CrossRoad, Revolution, Crisis,...



THE problem:
How to program these machines!!!!

Jesús Labarta, MP, 2008



Modelos de programación

■ Componentes

- Datos
- Procesos
- Comunicación
- Sincronización
- Entrada/salida

Jesús Labarta, MP, 2008



Memory

■ is THE problem

- 16 → 32 → 64 bits:
 - ✓ A simple thing for a computer architect
 - ✓ has caused huge pain for everybody !!!!

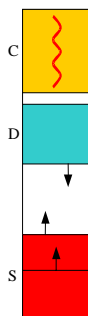
■ Some issues

- Addressability
- Association
- Consistency


Jesús Labarta, MP, 2008



Secuencial



Variable globales: 

Variables dinámicas(locales): 
Visibilidad local
Se puede pasar puntero al llamar funciones
Vida limitada

Librerías:

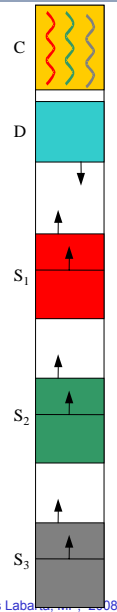
malloc: 
Visible global

Problemas:
Desbordamiento pilas
Errores difíciles detectar (efecto variable y lejano)




Jesús Labarta, MP, 2008



Memoria compartida: espacio @



Variable globales:
Unica copia 

Variables dinámicas(locales):
Acceso normal: privadas   
Accesibles globalmente (pasar puntero)

Librerías:
Exclusión mutua

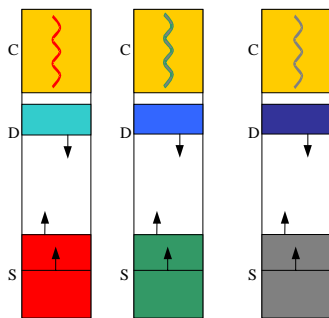
malloc:
Visible global 

Problemas:
Desbordamiento pilas
Errores difíciles detectar (efecto variable y lejano)

Jesús Labarta, MP, 2008



Memoria Distribuida: espacio @



Variables globales:
Replicadas: misma dirección lógica
mantener consistencia
Distribuidas: misma dirección, distinto objeto

Variables dinámicas(locales):
privadas (no posibilidad compartir)
~ mismas direcciones lógicas

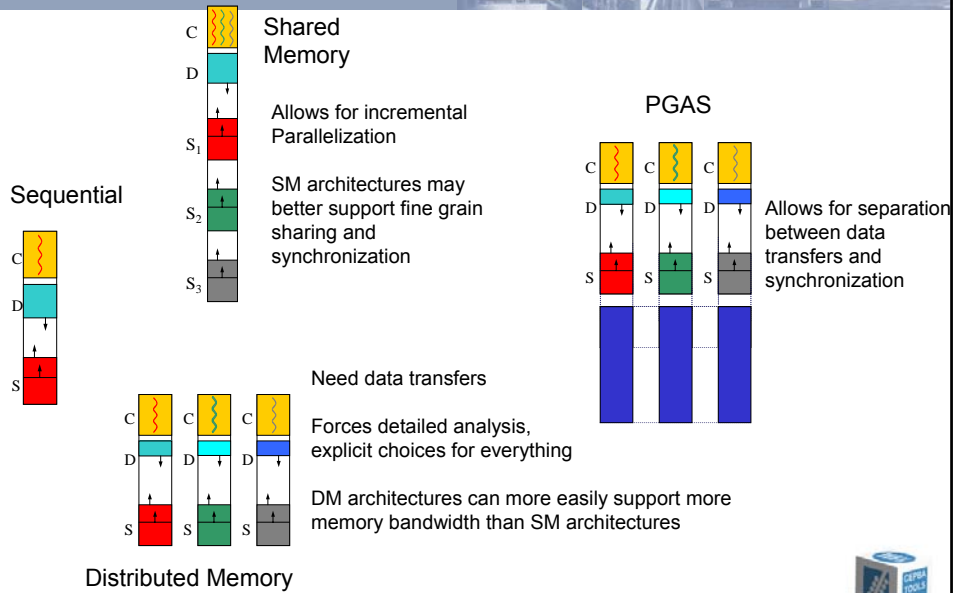
Librerías:
Locales.

malloc: local
~ mismas direcciones lógicas

Jesús Labarta, MP, 2008



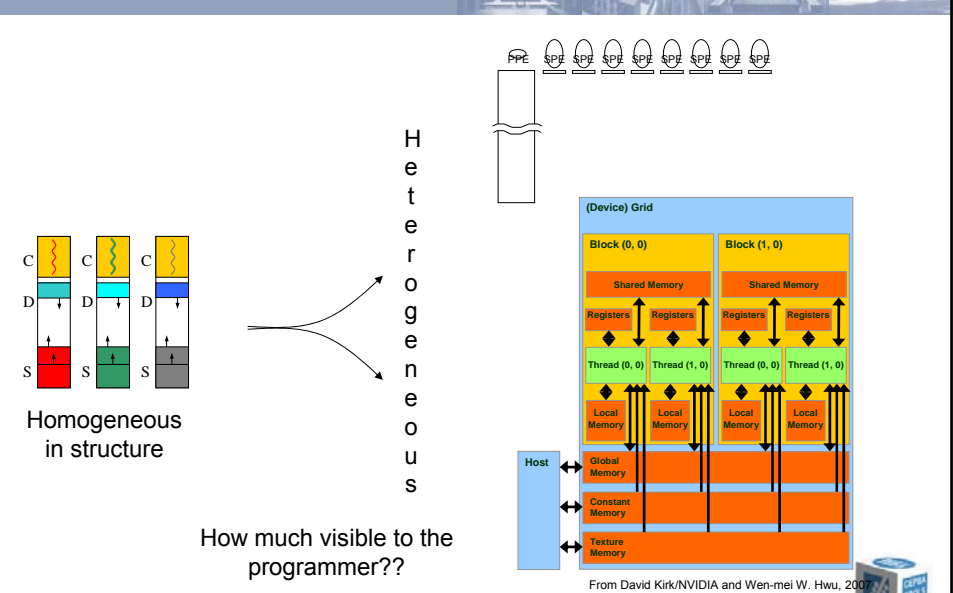
Address spaces in parallel programming



Jesús Labarta, MP, 2008



Address spaces in multicores



Jesús Labarta, MP, 2008



Modelos

■ Memoria compartida

- Pthreads
- OpenMP
 - ✓ <http://www.openmp.org>
 - ✓ <http://www.comunity.org>

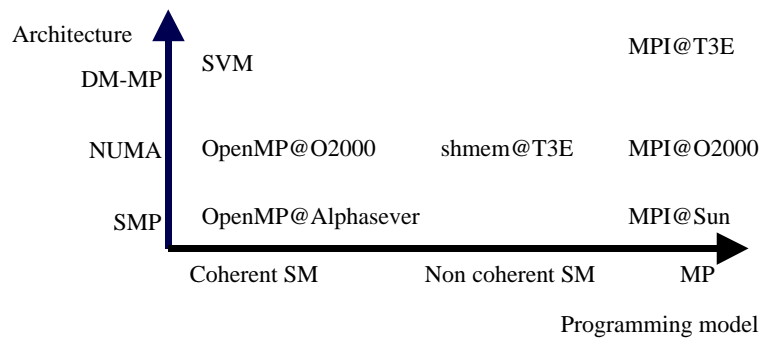
■ Memoria distribuida

- MPI
 - ✓ <http://www-unix.mcs.anl.gov/mpi/>

Jesús Labarta, MP, 2008



Espectro



Jesús Labarta, MP, 2008



Paralelización

■ Fases

- Descomposición
- Asignación
- Orquestación
- Mapeo

← Algorítmico, ideas

← Programación, sistema

Jesús Labarta, MP, 2008



■ Aspectos básicos de programación en los distintos modelos

Jesús Labarta, MP, 2008

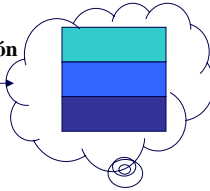


Pthreads: programación

Memoria compartida



Distribución
de cálculo



```
Double A(30,30)
integer i,j
```

```
Do j = 1,30
  Do i =1,30
    A(i,j) = 2* A(i,j)
    A(i,j) = j* A(i,j)
    A(i,j) = i* A(i,j)
```



```
Double A(30,30)
integer thid, j
```

```
main() {
  Do j = 1,30
    Do thid=1,num_threads
      iinf = f(thid)
      isup = f(thid)
      create_thread(Loop_body,iinf,isup)
    }
}
```

```
Loop_body(iinf,isup) {
  Do i =iinf,isup
    A(i,j) = 2* A(i,j)
    A(i,j) = j* A(i,j)
    A(i,j) = i* A(i,j)
  }
}
```

Jesús Labarta, MP, 2008

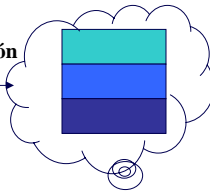


OpenMP: programación

Secuencial

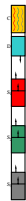


Distribución
de cálculo



```
Double A(30,30)
integer i,j
```

```
Do j = 1,30
  Do i =1,30
    A(i,j) = 2* A(i,j)
    A(i,j) = j* A(i,j)
    A(i,j) = i* A(i,j)
```



```
Double A(30,30)
integer i,j
```

```
Do j = 1,30
  C$OMP PARALLEL DO
  C$OMP PRIVATE (i)
  Do i =1,30
    A(i,j) = 2* A(i,j)
    A(i,j) = j* A(i,j)
    A(i,j) = i* A(i,j)
```

Jesús Labarta, MP, 2008

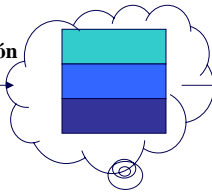


MPI: programación

Secuencial



Distribución de datos

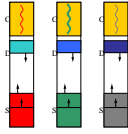


Memoria distribuida



```
Double A(30,30)
integer i,j
```

```
Do j = 1,30
  Do i =1,30
    A(i,j) = 2* A(i,j)
    A(i,j) = j* A(i,j)
    A(i,j) = i* A(i,j)
```



```
Double A(10,30)
```

```
myrank = quien_soy()
Do j = 1,30
  Do li =1,10
    i=li+myrank*10
    A(li,j) = 2* A(li,j)
    A(li,j) = j* A(li,j)
    A(li,j) = i* A(li,j)
```

Jesús Labarta, MP, 2008



Economía

■ Precio / rendimiento

■ Costes

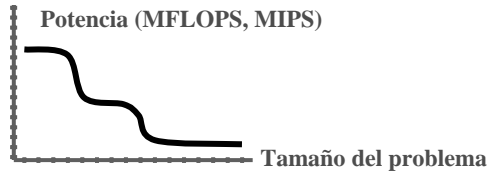
- Adquisición
- Mantenimiento
- Espacio y consumo
- Gestión
 - ✓ systems

Jesús Labarta, MP, 2008

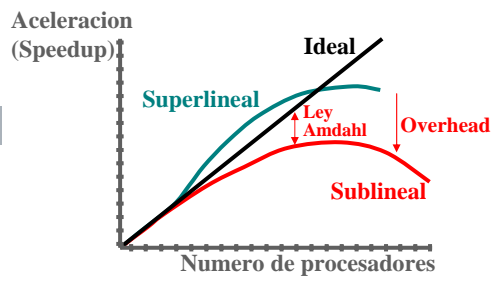


Que se puede esperar?

1 Procesador



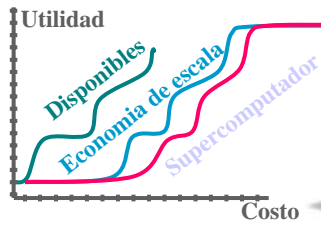
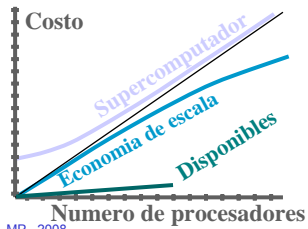
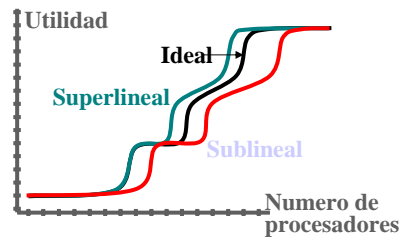
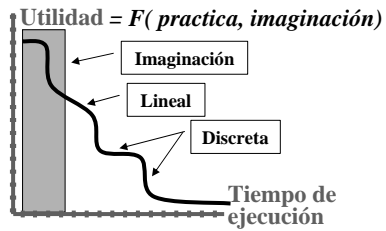
Paralelismo



Jesús Labarta, MP, 2008



Que se puede esperar?



Jesús Labarta, MP, 2008



Economía

■ Necesidad:

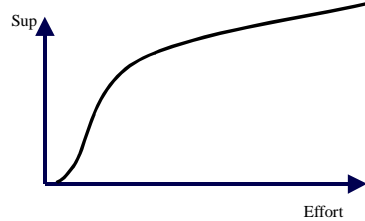
- Usuario \cong gas
 - ✓ $PV=nRT$
 - $P \propto 1/V \Rightarrow$ Solo si esta muy comprimido empuja de verdad
- Usuario \neq gas
 - Una vez expandido es difícil volver atrás
- Tendencias
 - ✓ COTS: procesadores, redes, software

Jesús Labarta, MP, 2008

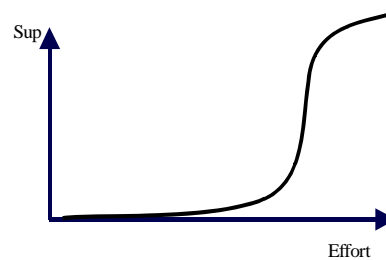


Economía: desarrollo

■ Memoria comparada



■ Memoria distribuida



Jesús Labarta, MP, 2008



Herramientas

■ Analizar una aplicación (sweep3d)

- Secuencial
- Paralela

■ Profile

- perfex
- ssrun

■ Traceo

- Paraver
- Dimemas
- Paradyn

Jesús Labarta, MP, 2008



Herramientas

■ Analizar rendimiento de aplicaciones

- Donde están sus cuellos de botella
- Información para optimizarla

■ Fases:

- Adquisición datos
- Análisis

Jesús Labarta, MP, 2008



Adquisición de datos

■ Instrumentación

- Qué:
 - ✓ código usuario
 - ✓ librería
 - ✓ sistema
- Quien:
 - ✓ manual
 - ✓ compilador
 - ✓ automático
- Método:
 - ✓ trazo / muestreo (correlado con ...)
 - ✓ Post-mortem / on-line

Jesús Labarta, MP, 2008



Adquisición de datos

■ A tener en cuenta

- Perturbación de la ejecución:
 - ✓ Causa: tiempo instrumentación, polución cache, almacenamiento,...
 - ✓ Minimización: implementación, soporte hardware
 - ✓ Efecto: conocerlo. Vivir con el.
- Tamaño datos:
 - ✓ Problema almacenamiento y postproceso
 - ✓ Control de la instrumentación: manual, automático
 - ✓ Procesado on-line
 - Overhead
 - Reducción información

Jesús Labarta, MP, 2008



Análisis

■ Objetivo

- Transmitir información al usuario: maximizar

■ Enfoque

- Percepción
- Cuantificación

■ Presentación

- Qué: magnitudes => vista
 - ✓ actividad temporal, fallos de cache, rutina, ...
- Cómo: modelo de presentación
 - ✓ Textual
 - ✓ Gráfico
 - ✓ sonoro,...

Jesús Labarta, MP, 2008



Aplicaciones

■ Solver iterativo

■ Radix sort

■ LU

Jesús Labarta, MP, 2008

