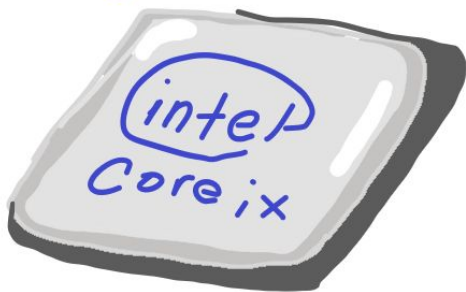


Curso Avanzado NLHPC



Ley de Moore



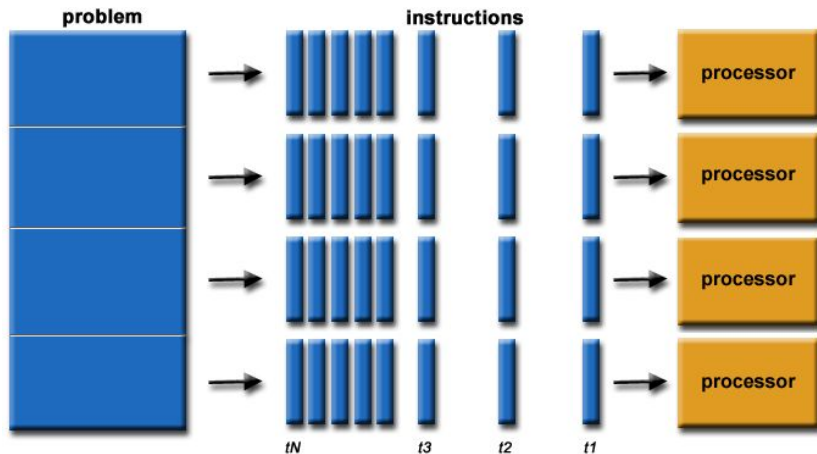
Aproximamente cada dos años se duplica
el numero de transistores en un
microprocesador.

GORDON MOORE, ES
CO- FUNDADOR DE
INTEL

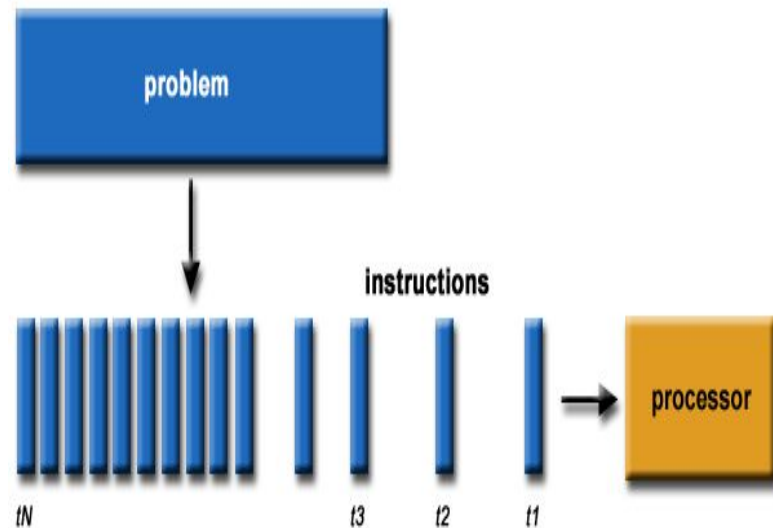
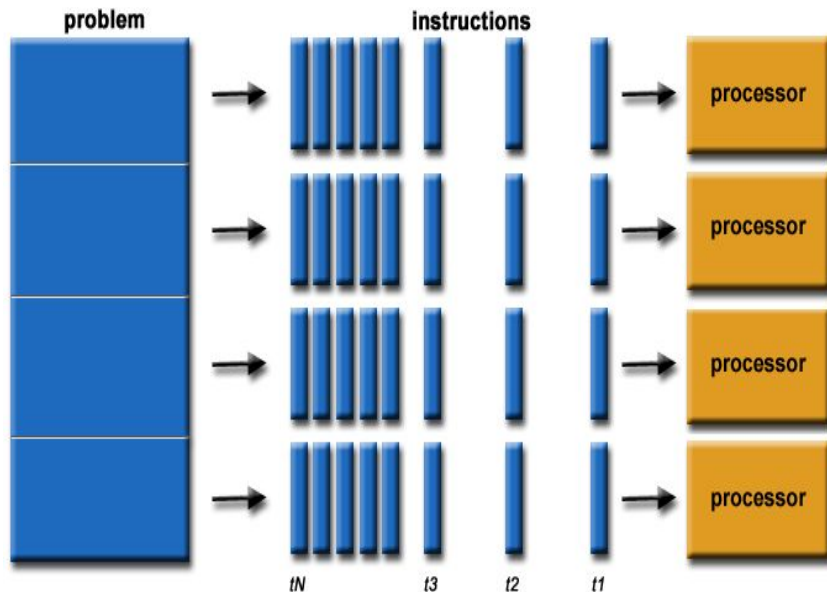


Computación Paralela

- Uso simultáneo de múltiples recursos computacionales
- Un problema se divide en partes discretas que se pueden resolver simultáneamente
- Cada parte se descompone en una serie de instrucciones
- Las instrucciones de cada parte se ejecutan simultáneamente en diferentes procesadores
- Se emplea un mecanismo global de control y coordinación

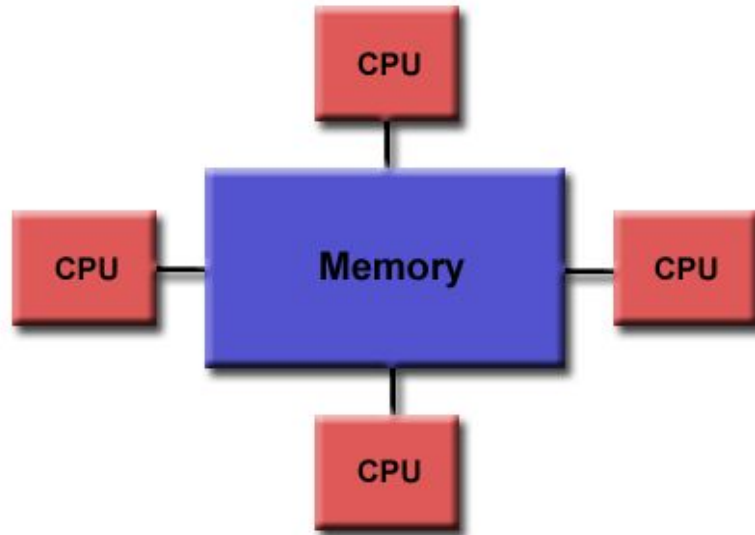


Computación Paralela Vs Computación Serial



Arquitecturas de memoria de computación Paralela: Memoria Compartida

- Los procesos comparten un espacio de memoria común
- Escriben y leen de manera asíncrona
- No es necesario especificar cómo se comunican los datos entre las tareas
- Se usan semáforos o locks para controlar el acceso a la memoria compartida



OpenMP

- Es una interfaz de programación de aplicaciones para la programación multiproceso de memoria compartida
- Permite añadir concurrencia: Operar actividades al mismo tiempo
- Se compone de:
 - Directivas de compilación
 - Rutinas de biblioteca
 - Variables de entorno
- Modelo de programación portable y escalable
- Proporciona a los desarrolladores una interfaz simple y flexible para el desarrollo de aplicaciones paralelas

Ejecución de Simulaciones en Slurm - Trabajos paralelos (OpenMP)

- Los trabajos paralelos diseñados para ejecutarse en un sistema multi-core (shared memory) requieren especificar el número de núcleos (-c 44) a utilizar:

```
#!/bin/bash
# -----SLURM Parameters-----
#SBATCH -J paralelo_openmp
#SBATCH -p general
#SBATCH -n 1
#SBATCH -c 44
#SBATCH --mem=192000
#SBATCH --mail-user=eguerre@nlhpc.cl
#SBATCH --mail-type=ALL
#SBATCH -o paralelo_openmp_%j.out
#SBATCH -e paralelo_openmp_%j.err

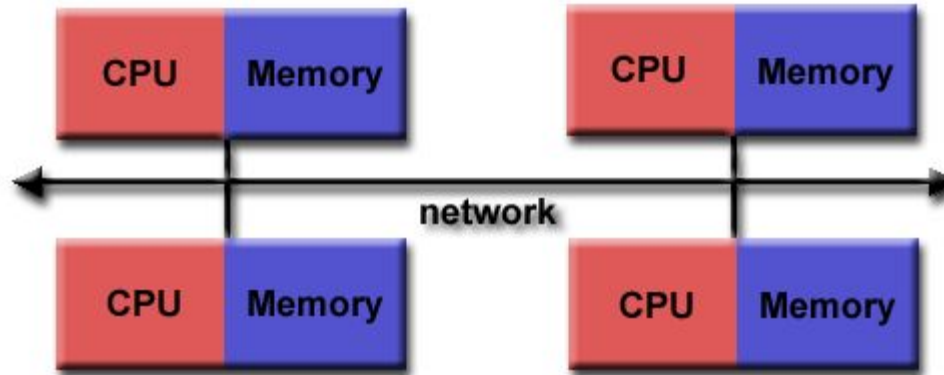
#-----Toolchain-----
ml purge
ml intel/2019b
# -----Módulos-----
ml Programa
# -----Comandos-----
./binario entrada
```

Ejercicio 1

1. Descarga el siguiente código: [pi_omp.c](#) en tu directorio de trabajo.
2. Teniendo en cuenta que es un código OpenMP, compílalo mediante el comando `icc pi_omp.c -o pi_omp -qopenmp`.
3. Crea un script de ejecución para lanzarlo con sbatch, con las siguientes consideraciones: a) La partición a lanzar es slims. b) Ejecuta 1 proceso. c) Este único proceso debe tener 20 hilos de ejecución. d) Supón que cada hilo ocupa 2400MB de memoria RAM. e) El binario a ejecutar es pi_openmp.
4. Anota el script que has usado.
5. Anota la salida de la ejecución.
6. Repite el ejercicio variando el número de hilos, ¿qué ocurre con el tiempo de ejecución?

Arquitecturas de memoria de computación Paralela: Memoria Distribuida

- Esta arquitectura se basa en múltiples procesadores con su propia memoria física privada
- Las tareas pueden operar solo con información local
- Se prefiere esta arquitectura ya que se pueden añadir múltiples unidades de procesamiento
- Requiere una red de comunicación para conectar la memoria entre procesadores
- Las tareas intercambian datos por medio del paso y recepción de mensajes



MPI

- Es una especificación para programación de pasos de mensajes
- Proporciona una librería de funciones para C, C++ o Fortran (Existen implementaciones para Python y R)
- Características:
 - Estandarización
 - Portabilidad (Multiprocesadores, multicomputadoras)
 - Buenas prestaciones
 - Múltiples implementaciones (MPICH, **OpenMPI**, **IMPI**, LAM-MPI, MVAPICH)
- El usuario escribirá su aplicación como un proceso secuencial del que se lanzarán varias instancias que cooperan entre sí

Ejecución de Simulaciones en Slurm - Trabajos paralelos (MPI)

- En el caso de trabajos en paralelo con MPI, es necesario especificar la cantidad de tareas en total (-n 132) y cuantas queremos que entren por nodo (--ntasks-per-node=44):

```
#!/bin/bash
# -----SLURM Parameters-----
#SBATCH -J paralelo_mpi
#SBATCH -p general
#SBATCH -n 132
#SBATCH --ntasks-per-node=44
#SBATCH --mem=192000
#SBATCH --mail-user=eguerre@nlhpc.cl
#SBATCH --mail-type=ALL
#SBATCH -o paralelo_mpi_%j.out
#SBATCH -e paralelo_mpi_%j.err

#-----Toolchain-----
ml purge
ml intel/2019b
# -----Módulos-----
ml WRF/4.1-dm+sm
# -----Comandos-----
srun wrf.exe|
```

Ejercicio 2

1. Descarga el siguiente código: [pi_mpi.c](#) en tu directorio de trabajo.
2. Teniendo en cuenta que es un código MPI, compílalo mediante el comando `mpicc pi_mpi.c -o pi_mpi`.
3. Crea un script de ejecución para lanzarlo con sbatch, con las siguientes consideraciones: a) La partición a lanzar es slims. b) Ejecuta 20 procesos. c) Deben entrar 10 procesos por cada nodo. d) Supón que cada tarea ocupa 2400MB de memoria RAM. e) El binario a ejecutar es `pi_mpi`.
4. Anota el script que has usado.
5. Anota la salida de la ejecución.
6. Repite el ejercicio variando el número de procesos, ¿qué ocurre con el tiempo de ejecución?

Ejecución de Simulaciones en Slurm - Trabajos paralelos Híbrido (MPI+OpenMP)

- En el caso de trabajos híbridos MPI+OpenMP, es necesario declarar el número de procesos MPI o tareas (-n 2) y además el número de cpus a asignar por cada una de esas tareas (-c 44):

```
#!/bin/bash
# -----SLURM Parameters-----
#SBATCH -J hibrido
#SBATCH -p general
#SBATCH -n 2
#SBATCH --ntasks-per-node=1
#SBATCH -c 44
#SBATCH --mem=192000
#SBATCH --mail-user=eguerre@nlhpc.cl
#SBATCH --mail-type=ALL
#SBATCH -o hibrido_%j.out
#SBATCH -e hibrido_%j.err

#-----Toolchain-----
ml purge
ml intel/2019b
# -----Módulos-----
ml Programa
# -----Comandos-----
srun ./programa
```

Ejercicio 3

1. Descarga el siguiente código: [hello_hybrid.c](#) en tu directorio de trabajo.
2. Teniendo en cuenta que es un código MPI-OpenMP híbrido, compílalo mediante el comando `mpiicc hello_hybrid.c -o hello_hybrid -qopenmp`.
3. Crea un script de ejecución para lanzarlo con sbatch, con las siguientes consideraciones: a) La partición a lanzar es slims. b) Ejecuta 2 procesos. c) Cada proceso debe de tener 20 hilos de ejecución. d) Supón que cada proceso ocupa 2400MB de memoria RAM. e) El binario a ejecutar es hello_hybrid.
4. Anota el script que has usado.
5. Anota la salida de la ejecución.

Ejecución de simulaciones en Slurm - Parámetros

Comando SLURM	Descripción
--mem-per-cpu= <megabytes>	Memoria requerida para el trabajo por CPU (en MegaBytes). El valor predeterminado es 1024 MB.
-p <partition>, --partition= <partition>	Enviar un trabajo a una partición específica.
-n, --ntasks=<cantidad de tareas>	Número de tareas que serán asignadas para el trabajo.
-c <cpus>	Esta es la cantidad de CPU que necesita su trabajo. Tenga en cuenta que SLURM es relativamente generoso con las CPU, y el valor especificado aquí es el número "mínimo" de CPU que se asignará a su trabajo. Si hay CPU adicionales disponibles en un nodo más allá de lo solicitado, su trabajo recibirá esas CPU hasta que otros trabajos las necesiten. El valor predeterminado es 1 CPU. Intentar usar más CPU de las que se le asignaron dará como resultado que sus procesos adicionales se turnen en la misma CPU (ralentizando su trabajo).
-J <name>, --jobname= <name>	Especifica un nombre a tu trabajo.
--mail- type=BEGIN,END,FAIL,ALL and --mail-user= <emailAddress>	Enviar por correo electrónico cuando su trabajo comienza / termina / falla. Puede especificar varios valores para esto (separados por comas) si es necesario.
-o <STDOUT_log>, -- output=<STDOUT_log>	Redirija la salida a los archivos de registro que especifique. Por defecto, ambos, STDOUT and STDERR son enviados a este archivo. Puedes especificar %j como parte del nombre de archivo de registro para indicar la ID del trabajo (como ejemplo, "#SBATCH -o output_%j.o" redirigirá la salida a "output_123456.o").
-e <STDERR_log>, --error= <STDERR_log>	Redirigir STDERR a un archivo separado. Funciona exactamente igual que "-o".
-t <days- hours:minutes:seconds>	Walltime para tu trabajo. La duración del Walltime es el tiempo que espera que su trabajo se ejecute.
-a, --array=<índices>	Envía una lista (arreglo) de trabajos idénticos. Solo aplica para sbatch.

Ejecución de Simulaciones en Slurm - Trabajos secuenciales

- En este ejemplo consideramos el caso de enviar una tarea utilizando 1 core:

```
#!/bin/bash
# -----SLURM Parameters-----
#SBATCH -J secuencial
#SBATCH -p slims
#SBATCH -n 1
#SBATCH --ntasks-per-node=1
#SBATCH -c 1
#SBATCH --mem=48000
#SBATCH --mail-user=eguerra@nlhpc.cl
#SBATCH --mail-type=ALL
#SBATCH -o secuencial_%j.out
#SBATCH -e secuencial_%j.err

#-----Toolchain-----
ml purge
ml intel/2019b
# -----Módulos-----
ml Python/3.8.2
# -----Comandos-----
python3 programa.py
```


Ejecución de Simulaciones en Slurm - Trabajos secuenciales (arrays)

- Un job array en Slurm es una colección de trabajos idénticos que sólo difieren entre sí por un parámetro. La dimensión del arreglo se define con el parámetro `--array=1-n`, donde `n` es la cantidad de jobs. Se puede definir la cantidad de procesos simultáneos de la siguiente forma: `1-100%10`. A continuación un ejemplo:

```
#!/bin/bash
# -----SLURM Parameters-----
#SBATCH -J secuencial_array
#SBATCH -p slims
#SBATCH -n 1
#SBATCH --ntasks-per-node=1
#SBATCH -c 1
#SBATCH --mem=48000
#SBATCH --mail-user=eguerra@nlhpc.cl
#SBATCH --mail-type=ALL
#SBATCH --array=1-100
#SBATCH -o secuencial_array_%A_%a.out
#SBATCH -e secuencial_array_%A_%a.err

#-----Toolchain-----
ml purge
ml intel/2019b
# -----Módulos-----
ml Python/3.8.2
# -----Comandos-----
python3 programa.py $SLURM_ARRAY_TASK_ID
```

Ejercicio 4

1. Descarga el siguiente código: [average.py](#) en tu directorio de trabajo.
2. Crea un script de ejecución para lanzarlo con sbatch, con las siguientes consideraciones: a) La partición a lanzar es slims. b) Ejecuta 1 tarea. c) Esta única tarea debe tener asignados 2 core. d) Supón que cada tarea ocupa 2400MB de memoria RAM. e) Debes utilizar la versión de Python más reciente. f) La dimensión del arreglo es de 100 y deben entrar de a 10 simultáneos
4. Anota el script que has usado.
5. Anota la salida de la ejecución.

Ejecución de Simulaciones en Slurm - Trabajos paralelos (GPU)

- Los trabajos que utilizarán las GPUs deben indicar la cantidad a utilizar (`--gres=gpu:1`)
Cada nodo tiene 2 GPUs:

```
#!/bin/bash
# -----SLURM Parameters-----
#SBATCH -J GPU
#SBATCH -p gpus
#SBATCH -n 44
#SBATCH --gres=gpu:1
#SBATCH --ntasks-per-node=44
#SBATCH -c 1
#SBATCH --mem=192000
#SBATCH --mail-user=eguerrea@nlhpc.cl
#SBATCH --mail-type=ALL
#SBATCH -o GPU_%j.out
#SBATCH -e GPU_%j.err

#-----Toolchain-----
ml purge
ml fosscuda/2019b
# -----Módulos-----
ml NAMD/2.13-mpi
# -----Comandos-----
srun namd2|
```

Ejercicio 5

1. Descarga los siguientes archivos en tu directorio de trabajo: [mulBy2.cu](#), [Makefile](#)
2. Teniendo en cuenta que es un código en cuda y que tenemos un Makefile, compílalo mediante el comando **make** en tu directorio de trabajo
3. Crea un script de ejecución para lanzarlo con sbatch, con las siguientes consideraciones: a) La partición a lanzar es gpus. b) Ejecuta 1 proceso. c) Cada proceso debe de tener 1 core. d) Supón que cada proceso ocupa 2400MB de memoria RAM. e) El binario a ejecutar es mulBy2.
4. Anota el script que has usado.
5. Anota la salida de la ejecución.

Dependencias de trabajos

Las dependencias de trabajos se utilizan para aplazar el inicio de un trabajo hasta que se satisfagan las dependencias especificadas. Se especifican con la opción `--dependency` en el siguiente formato:

```
sbatch --dependency=<type:job_id[:job_id][,type:job_id[:job_id]]> ...
```

Los tipos de dependencias son las siguientes:

`after:jobid[:jobid...]`

el trabajo puede empezar después de que los trabajos especificados comiencen

`afterany:jobid[:jobid...]`

el trabajo puede empezar después de que los trabajos especificados terminen

`afternotok:jobid[:jobid...]`

el trabajo puede empezar después que los trabajos especificados terminan fallidamente

`afterok:jobid[:jobid...]`

el trabajo puede empezar después que los trabajos especificados terminan exitosamente

Ejemplo de dependencias de trabajos - Ejemplo

La manera más simple de usar una dependencia del tipo afterok:

```
[prueba@leftrar1 ~]$ sbatch job1.sh
```

Submitted batch job 21363626

```
[prueba@leftrar1 ~]$ sbatch --dependency=afterok:21363626 job2.sh
```

Ahora cuando job1.sh termine correctamente, el job2.sh entrará en ejecución. Si job1.sh termina fallidamente, job2.sh no entrará en ejecución nunca pero sí quedará en cola (debe cancelarse manualmente el trabajo).

Programación de tareas usando crontab

Crontab es un planificador de tareas donde se guarda un listado de comandos a ejecutar en un tiempo determinado por el usuario.

Para acceder a crontab utilice el siguiente comando:

```
[prueba@leftrar1 ~]$ crontab -e
```

Esto le permitirá editar su archivo crontab, la estructura base es la siguiente:

```
# Example of job definition:
# .----- minute (0 - 59)
# | .----- hour (0 - 23)
# | | .----- day of month (1 - 31)
# | | | .----- month (1 - 12) OR jan,feb,mar,apr ...
# | | | | .---- day of week (0 - 6) (Sunday=0 or 7) OR sun,mon,tue,wed,thu,fri,sat
# | | | | |
# * * * * * command to be executed
```

Programación de tareas usando crontab - Ejemplo

Primero debemos especificar el tiempo, en este caso la tarea será ejecutada todos los días cada 20 minutos, luego debemos obtener el archivo `/etc/profile` que contiene las definiciones comunes para los usuarios, posterior a ello, declaramos el comando `sbatch` para iniciar nuestra tarea con el gestor de recursos `slurm`, para finalizar, indicamos el script enviado por el usuario y guardamos los cambios en el archivo:

```
*/20 * * * * . /etc/profile ; sbatch /home/prueba/ejemplo/script.sh
```

Para revisar el listado de tareas existentes en nuestro crontab, ejecutamos el comando:

```
[prueba@leftraru1 ~]$ crontab -l
```

```
*/20 * * * * . /etc/profile ; sbatch /home/prueba/ejemplo/script.sh
```

Para borrar el contenido de nuestro crontab:

```
[prueba@leftraru1 ~]$ crontab -r
```


Monitoreo de Simulaciones - htop

- Puede ingresar a través de ssh a un nodo en donde tenga una tarea en ejecución y ejecutar **htop**:

```
1 [|||||100.0%] 6 [ 0.0%] 11 [ 0.0%] 16 [ 0.0%]
2 [ 0.0%] 7 [ 0.0%] 12 [ 0.0%] 17 [ 0.0%]
3 [|||||100.0%] 8 [ 0.0%] 13 [ 0.0%] 18 [ 0.0%]
4 [ 0.0%] 9 [ 0.7%] 14 [ 0.0%] 19 [ 0.0%]
5 [ 0.0%] 10 [ 0.0%] 15 [ 0.0%] 20 [ 0.0%]

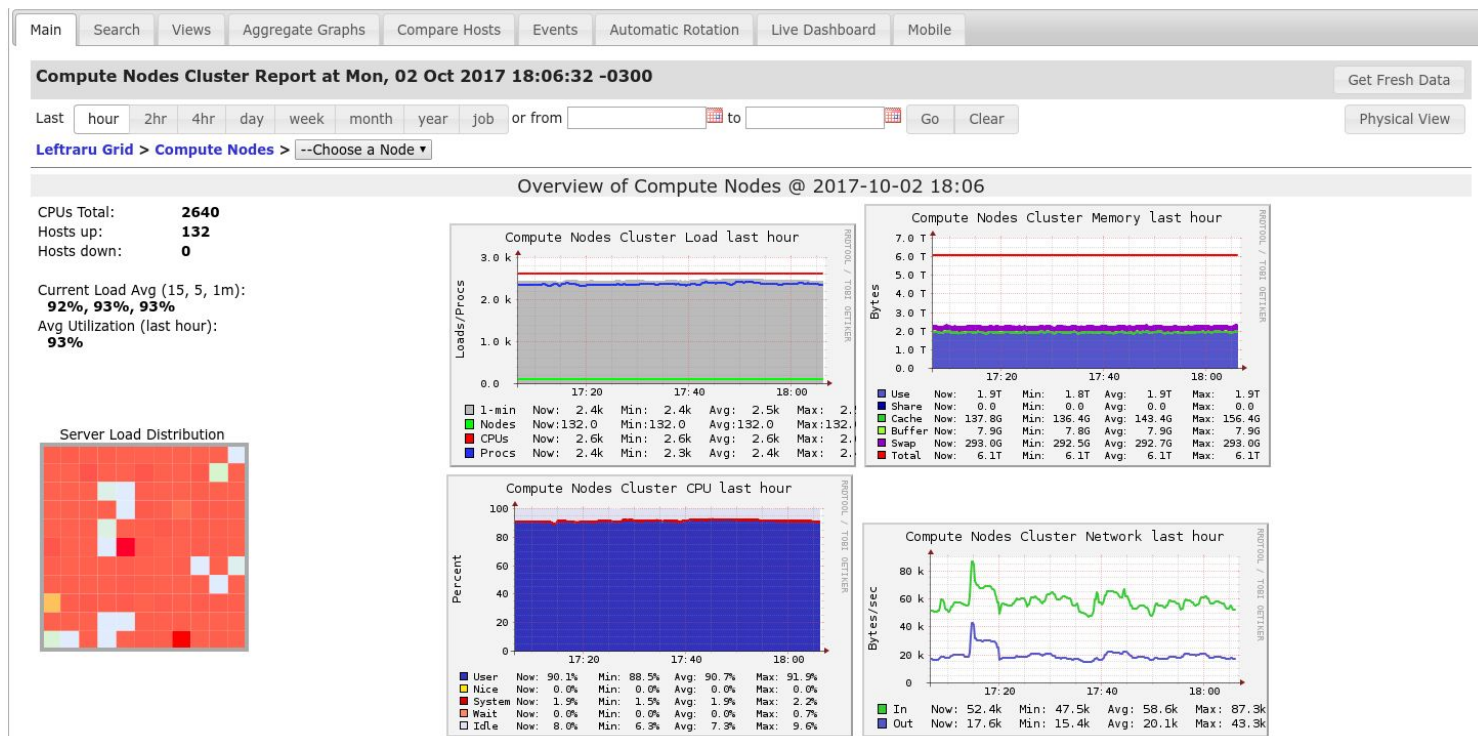
Mem[||||| 1.56G/62.7G] Tasks: 44, 39 thr: 3 running
Swp[ 0K/62.5G] Load average: 2.00 2.01 2.05
Uptime: 4 days, 00:36:11

PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
18305 nperinet 20 0 19020 3768 892 R 100.0 0.0 2h52:34 ./LL RK4.x
18335 nperinet 20 0 17104 1584 892 R 100.0 0.0 2h47:30 ./LL RK4.x
18605 root 20 0 121M 2432 1476 R 0.0 0.0 0:00.05 htop
1 root 20 0 185M 5068 2384 S 0.0 0.0 0:11.06 /usr/lib/systemd/systemd --switched-root --system --de
669 root 20 0 112M 2000 1564 S 0.0 0.0 0:00.00 /bin/bash
671 root 20 0 36816 7236 6908 S 0.0 0.0 0:00.78 /usr/lib/systemd/systemd-journald
726 root 20 0 43808 2428 1268 S 0.0 0.0 0:00.83 /usr/lib/systemd/systemd-udev
1012 root 16 -4 51188 1620 1236 S 0.0 0.0 0:00.05 /sbin/auditd -n
1002 root 16 -4 51188 1620 1236 S 0.0 0.0 0:00.25 /sbin/auditd -n
1206 root 20 0 448M 10956 6968 S 0.0 0.0 0:00.00 /usr/sbin/NetworkManager --no-daemon
1209 root 20 0 448M 10956 6968 S 0.0 0.0 0:00.09 /usr/sbin/NetworkManager --no-daemon
1139 root 20 0 448M 10956 6968 S 0.0 0.0 0:02.34 /usr/sbin/NetworkManager --no-daemon
1144 avahi 20 0 30220 1564 1300 S 0.0 0.0 0:00.44 avahi-daemon: running [cnf004.local]
1148 dbus 20 0 28824 1772 1352 S 0.0 0.0 0:00.44 /bin/dbus-daemon --system --address=systemd: --nofork
1166 root 20 0 198M 1236 776 S 0.0 0.0 0:00.00 /usr/sbin/gssproxy -D
1167 root 20 0 198M 1236 776 S 0.0 0.0 0:00.00 /usr/sbin/gssproxy -D
1168 root 20 0 198M 1236 776 S 0.0 0.0 0:00.00 /usr/sbin/gssproxy -D
1169 root 20 0 198M 1236 776 S 0.0 0.0 0:00.00 /usr/sbin/gssproxy -D
1170 root 20 0 198M 1236 776 S 0.0 0.0 0:00.00 /usr/sbin/gssproxy -D
1164 root 20 0 198M 1236 776 S 0.0 0.0 0:00.30 /usr/sbin/gssproxy -D

F1Help F2Setup F3Search F4Filter F5Free F6SortBy F7Nice F8Nice F9Kill F10Quit
```

Monitoreo de simulaciones - Ganglia

- Para monitorear nuestras tareas en ganglia, deben ir a <http://monitor.nlhpc.cl/ganglia/>



Monitoreo de Simulaciones - Ganglia

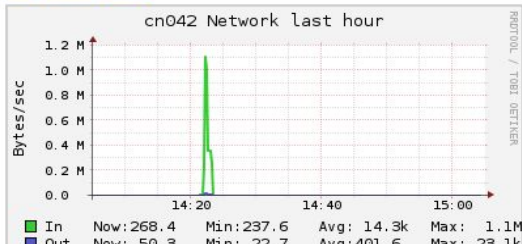
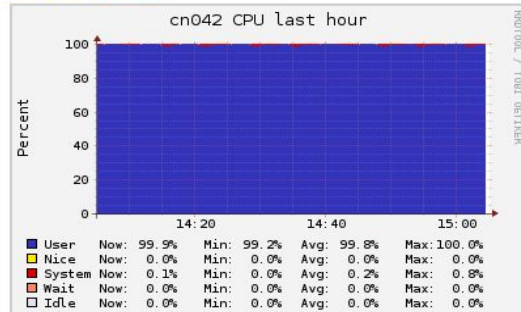
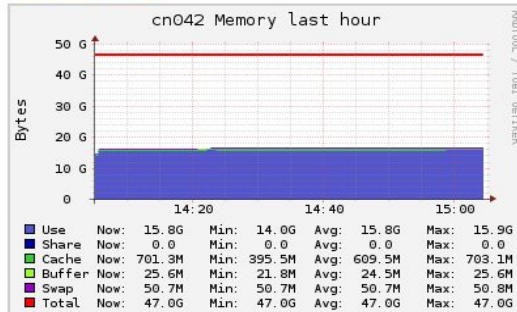
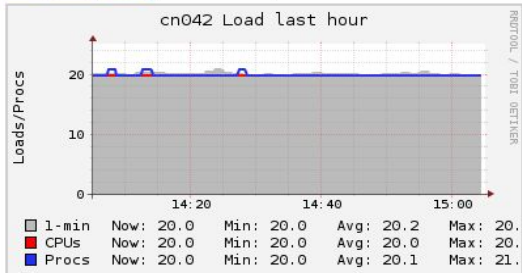
cn042 Host Report at Tue, 16 Aug 2016 15:04:59 -0300

Get Fresh Data

Last or from to

Leftraru Grid > Compute Nodes > cn042

Host Overview



Ejercicio 6

1. ¿Cuánto espacio de almacenamiento tienes ocupado en tu cuenta de usuario? ¿Qué comando debes usar para obtener dicha información?
2. Si quieres lanzar un programa en MPI en la infraestructura del NLHPC y tienes 120 cores disponibles, ¿cuánto deberías de reservar/usar?

Instalación y Compilación de aplicaciones - Compiladores

- La compilación es el proceso por el cual se traducen las instrucciones escritas en un determinado lenguaje de programación a lenguaje máquina
- Un compilador es un programa de computadora que “traduce” las instrucciones o código fuente desde un lenguaje a otro (usualmente código máquina)
- En ambientes HPC, se busca sacar el máximo rendimiento al hardware con las aplicaciones, por lo que el 99% de los programas se encuentra compilado
- En nuestro caso, disponemos de los siguientes set de herramientas de compilación:
 - **Intel Toolchain**
 - **GNU Toolchain**
 - Estos 2 set de herramientas de compilación permiten compilar programas escritos en C, C++, Fortran
- Recomendamos siempre hacer uso de las herramientas de Intel, ya que se encuentran optimizadas para utilizar las características de sus microprocesadores en toda su extensión

Instalación y Compilación de aplicaciones - Flags

- Los flags de optimización se usan para:
 - Disminuir la cantidad de mensajes de depuración
 - Aumentar los niveles de aviso de errores
 - Optimizar el código producido
- Algunos de los flags de optimización utilizados en nuestras aplicaciones:
 - **-O3**: Habilitar la optimización del compilador, donde 3 representa el nivel (máximo)
 - **-ipo**: La optimización interprocedural le permite al compilador analizar el código para determinar donde este puede beneficiarse de optimizaciones específicas
 - **-no-prec-div**: Mejora la precisión de las divisiones de punto flotante
 - **-static-intel**: Hace que las bibliotecas proporcionadas por intel se vinculen estáticamente
 - **-fPIC**: Determina si el compilador genera código independiente de la posición
 - **-qopenmp**: Permite que el paralelizador genere código multiproceso basado en directivas OpenMP

Instalación y Compilación de aplicaciones - Ejemplos de compilación

- La compilación más básica de un programa en C sería la siguiente:
 - **icc hello_world.c -o hello_world**
 - Esta compilación generará un ejecutable llamado “hello_world” el cual puede ser ejecutado de la siguiente forma: **./hello_world**
- Si queremos compilar un programa en C que soporte MPI, la compilación sería la siguiente:
 - **mpiicc hello_world.c -o hello_world -O3 -axCORE-AVX512,AVX,SSE4.2 -ipo -no-prec-div -static-intel -fPIC -qopenmp**
- La opción **-axCORE-AVX512,AVX,SSE4.2** la utilizamos para que nuestro ejecutable pueda ejecutarse en los distintos tipos de nodos que tenemos actualmente
- Compiladores disponibles: **gcc/icc/mpiicc** (C) **g++/icpc/mpiicpc** (C++), **gfortran/ifort/mpiifort** (Fortran)

Instalación de módulos en Python y R

- En el caso de que queramos instalar módulos de Python localmente en nuestro directorio utilizaremos el siguiente comando:
 - Cargar el módulo de Python: **m1 Python/3.8.3**
 - **pip install programa --no-binary :all: --user**
 - El parámetro **--no-binary :all:** intentará compilar el módulo a instalar
 - El parámetro **--user** instalará en nuestro directorio el módulo
- Para instalar paquetes de R localmente en nuestro directorio:
 - Cargamos el módulo de R: **m1 R/4.0.0**
 - **R CMD INSTALL -l /home/eguerra/R/library programa**
 - Para decirle a R que utilice los paquetes instalados por mí: **library("paquete", lib.loc="/home/eguerra/R/library")**

Problemas Frecuentes - Exceso de memoria

- Debemos tener en cuenta que la memoria RAM que SLURM reserva por defecto es de 1000MB (1GB). Un típico error de cancelación de tareas es por utilizar más de la memoria reservada:

```
/tmp/slurmd/job136839939/slurm_script: line 15: 23547 Killed                  ./programa.sh
slurmstepd: error: Detected 1 oom-kill event(s) in step 136839939.batch cgroup. Some of your processes
may have been killed by the cgroup out-of-memory handler.
```

- Si su tarea ocupa más de la memoria por defecto (1GB), se puede utilizar el siguiente parámetro en SLURM:
 - **#SBATCH --mem=2400**
- Esto hará que SLURM reserve 2400MB (2.4GB) por la totalidad del trabajo
- Para más detalles de nuestra infraestructura visitar:
http://wiki.nlhpc.cl/index.php/Recursos_Computacionales_NLHPC

Problemas Frecuentes - Subutilización de CPU y RAM

- Nuestro objetivo en el cluster es velar por el uso óptimo de los recursos, por lo que estamos continuamente monitoreando el uso de los nodos. En el caso de la subutilización de CPU, la condición para que una tarea sea cancelada es que la mitad del total de cores reservados muestre un porcentaje de uso menor o igual a un 30% en un lapso de 4 horas. Para el caso de la memoria RAM, lo que se busca es garantizar el uso de al menos un 70% del recurso reservado.
- El procedimiento en ambos casos es el siguiente:
 - Se envía una notificación a las 2 horas de ejecución de la tarea que subutilice recursos
 - Si mantiene el comportamiento en las próximas 2 horas, se envía una notificación informando la cancelación de la tarea
 - En el correo de cancelación se adjuntan los gráficos del uso de CPU y Memoria por cada nodo con el fin de que puedan realizar los cambios respectivos en las próximas tareas

Problemas Frecuentes - Ejemplo subutilización de CPU

- Tenemos la tarea de un usuario que reserva 4 cpu y 4800MB de memoria. La notificación recibida es la siguiente:

Estimada/o Pablo,

Durante los procesos rutinarios de monitoreo hemos observado el comportamiento de la tarea **18197288 (test_general)**, iniciada por el usuario **pflores**.

Actualmente presenta un comportamiento que podemos clasificar como: subutilización de recursos, dado que 4 de sus cores presentan menos del 2% de utilización en las últimas 2 horas.

Los datos recopilados son:

JobID	Start	End	Elapsed	Partition	NodeList	NCPUS	ReqMem
18197288	2020-05-24T15:03:24	Unknown	02:44:08	slims	cn042	4	4800Mn

Consideramos esta información relevante, ya que su tarea se comporta en forma poco eficiente, al no utilizar óptimamente los recursos reservados dentro del clúster.

- Si este comportamiento se repite en las siguientes 2 horas, la tarea será cancelada
- Al revisar el archivo report.log de la notificación, tenemos lo siguiente:

```
* cpu= 9 avg=0.75 std=0.43 / all avg=0.79 std=0.41
* cpu= 8 avg=1.08 std=0.40 / all avg=1.07 std=0.36
* cpu= 7 avg=1.08 std=0.40 / all avg=1.07 std=0.36
* cpu= 6 avg=1.00 std=0.00 / all avg=1.00 std=0.00
```

- Debido a que el uso de CPU es tan bajo, la recomendación es:
 - Reservar siempre lo realmente utilizado, en este caso 1 cpu

Problemas Frecuentes - Ejemplo subutilización de Memoria

- Tenemos la tarea de un usuario que reserva 12 cpu y 2.5GB por cada cpu, osea un total de 30GB:

Estimada/o Eugenio,

Durante los procesos rutinarios de monitoreo hemos observado el comportamiento de la tarea **18202593 (ORCA-a-bpa-e)**, iniciada por el usuario **eguerre**.

Actualmente presenta un comportamiento que podemos clasificar como: subutilización de recursos, dado que desde el inicio de la ejecución presenta menos del **70%** de utilización de memoria.

Los datos recopilados son:

JobID	Start	End	Elapsed	Partition	NodeList	NCPUS	ReqMem
18202593	2020-05-25T14:50:13	Unknown	02:10:47	slims	cn075	12	2.50Gc

- Si este comportamiento se repite en las siguientes 2 horas, la tarea será cancelada
- Al revisar el archivo report.log de la notificación, tenemos lo siguiente:

```
* mem max:16861 resv:30720
```

- Debido a que el uso de memoria sólo llega a 16.8GB:
 - Se debe reservar la memoria que realmente se utiliza:
 - **#SBATCH --mem=18000**

Gracias!

<https://www.nlhpc.cl/>