

Computación de Altas Prestaciones Seminario sobre GPGPUs



Contenidos

Sesión 1 Teoría: Introducción a Computación en GPUs con CUDA.

Sesión 2: Compilación, conceptos básicos

Sesión 3: Programación de algoritmos “trivialmente paralelos”

Sesión 4: Uso de la memoria “Shared”.
“Reducciones” en GPUs

Sesión 5: Optimización, temas avanzados, librerías

Tarjeta gráfica y su GPU

El objetivo fundamental de la tarjeta gráfica es encargarse de la visualización. Toman mucha fuerza al surgir los sistemas operativos gráficos (Windows, Mac)

- Visualización 2D
- Ejecución de video
- Gráficos 3D

Tarjeta gráfica y su GPU

- A finales de los 90 y comienzos de los 2000, las GPUs se van haciendo cada vez más potentes.
- Operaciones principales : rendering, shading,... operaciones sencillas que hay que realizar sobre muchos pixeles.
- Se desarrollan GPUs con muchos elementos de computación, programables, relativamente sencillos (sin ejecución out of order, sin ejecución especulativa, sin caches).
- Hasta ese momento, la CPU manda ordenes a GPU, que ejecuta ordenes gráficas. GPGPU surge cuando aparece la posibilidad de mandar resultados de vuelta a CPU

Tarjeta gráfica y su GPU

- La potencia de cálculo acumulada de todos esos cores empieza a ser mayor que la de las CPUs

Surge la idea de usarla para computación de carácter general.

Inicialmente:

- difíciles de programar (OpenGL, DirectX), poca flexibilidad (programar otros cálculos como si fueran operaciones de gráficos)

- poca precisión (no había números reales)

- Sin medios de depurar

CUDA

CUDA (*Compute Unified Device Architecture*) es (en el momento de presentación, 2006) una nueva arquitectura para tarjetas gráficas con elementos diseñados para cálculo general con GPUs.

-Para facilitar el acceso a esas tarjetas, Nvidia creó CUDA-C: Una serie de “añadidos” a C estándar para permitir el acceso

CUDA:

Definición actual de CUDA:

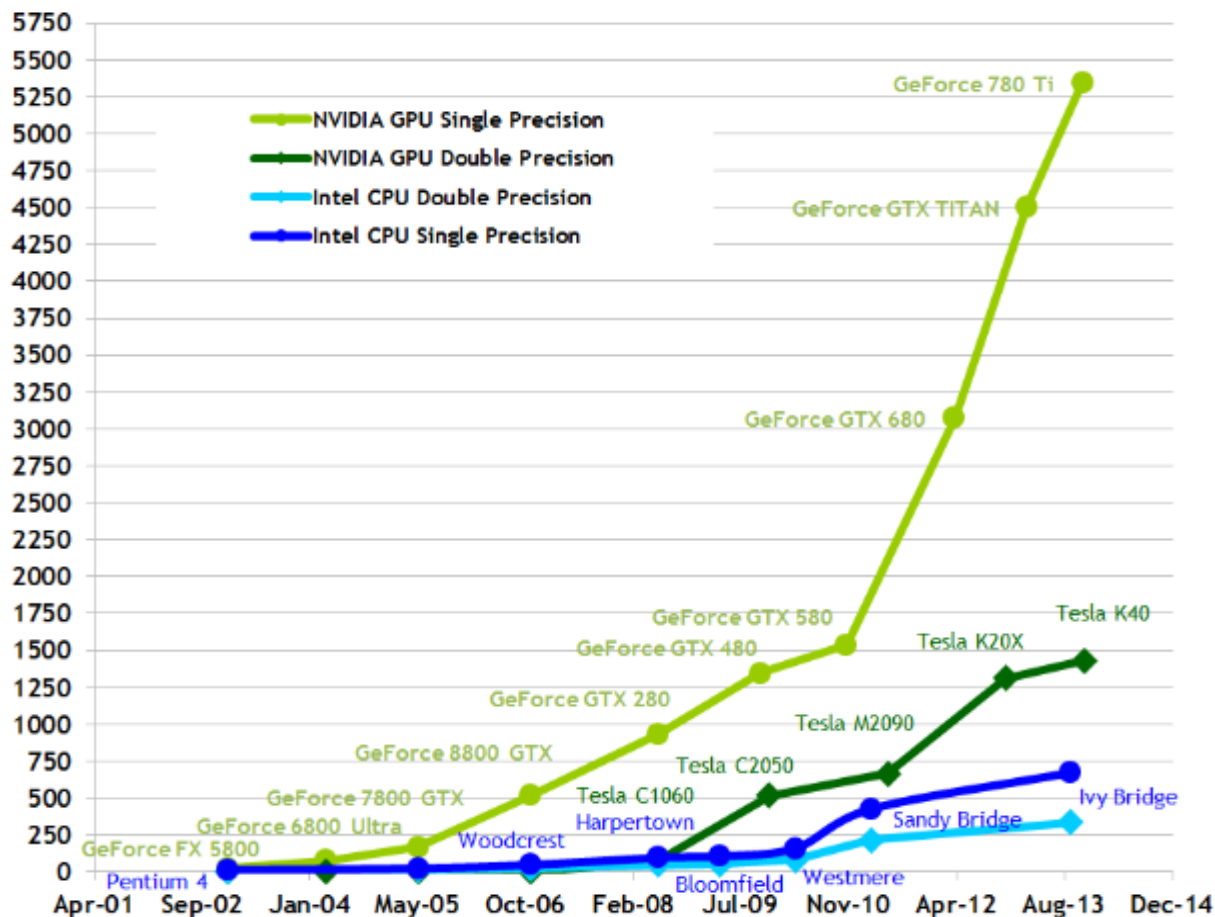
"a general purpose parallel computing platform and programming model that leverages the parallel compute engine in NVIDIA GPUs to solve many complex computational problems in a more efficient way than on a CPU."

-La principal forma de programar GPUs de Nvidia es mediante CUDA-C. Sin embargo, hay otros lenguajes que lo permiten: FORTRAN, DirectCompute, OpenACC, o desde Matlab.

Porque CUDA?

Figure 1. Floating-Point Operations per Second for the CPU and GPU

Theoretical GFLOP/s



<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#introduction>

Aplicaciones de CUDA

Se han desarrollado aplicaciones de CUDA prácticamente en cualquier campo en el que sea necesaria computación masiva:

- Medicina
- Biología
- Astrofísica
- Dinámica de fluidos computacional
- Estudio del clima
- Geofísica
- Economía
- Big Data
- Estadística
- ...

La lista no acaba, porque crece cada día.

Hay muchos casos de librerías aceleradas con CUDA, para cuyo uso no es necesario usar CUDA en absoluto.

Lista actualizada en:

<https://developer.nvidia.com/gpu-accelerated-libraries#signal>

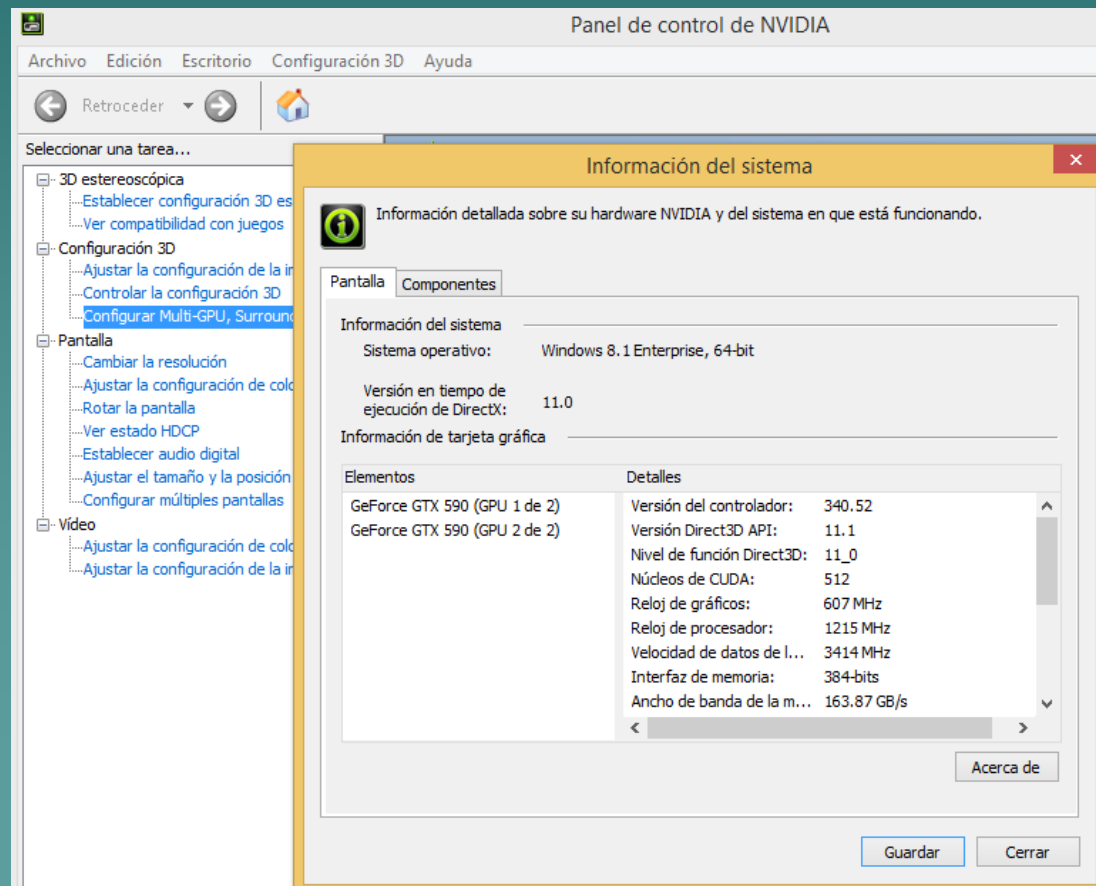
Alternativas a CUDA

OpenCL es un lenguaje bastante similar a CUDA. OpenCL es Open Source (apoyado sobre todo por AMD), mientras que CUDA es propiedad de Nvidia.

En general se considera que CUDA es una tecnología mas madura (y eficiente) que OpenCL. Por ejemplo, Matlab (a través del ParallelComputing Toolbox) soporta tarjetas gráficas que soporten CUDA. Las tarjetas Nvidia soportan también OpenCL.

Como podemos usar CUDA en nuestro ordenador?

1) Necesitamos una GPU NVIDIA compatible con CUDA
(Prácticamente todas las tarjetas Nvidia son ya compatibles con CUDA: GeForce, Quadro, Tesla...)



Como podemos usar CUDA en nuestro ordenador?

- 2) Tenemos que instalar un “device driver” apropiado (el último disponible para la tarjeta servirá).
Lo ideal es tener un sistema con una GPU para visualización y otra para cálculo, pero puede ser difícil de configurar.
- 3) Tenemos que descargar e instalar el “CUDA Toolkit” de la web de Nvidia.
- 4) Compilador de C compatible:
 - Windows: necesitamos tener instalada alguna versión apropiada de Visual Studio
 - En Linux no se necesita software “propietario”; hay que instalarlo en modo “root”

Podéis consultar “Cuda Quick Start Guide” en [poliformat](#)
NO podemos depurar programas CUDA con una GPU que se esté usando para controlar una interfaz gráfica de usuario.

Máquinas disponibles en el DSIC

1) PCS de laboratorios 0,6,7,8

2) Máquinas `knight.dsic.upv.es` y `gpu.dsic.upv.es`, a la que podeis acceder con ssh desde el dominio `dsic.upv.es`

-`knight.dsic.upv.es`: Intel Xeon con 24 cores + 1 Teslas K20c
Tesla K20c: 5Gb memoria, 13 Multiprocesadores con 192 cores
cada uno: 2496 cores

-`gpu.dsic.upv.es` : Core i9 (12 cores) + 2 Quadro RTX 5000 con 16 Gb y 3000 cores cada una

Para averiguar cuantas gpus tenemos en nuestro sistema, podemos ejecutar la aplicación `deviceQuery`(parte del conjunto de "Samples" que viene con el "CUDA Toolkit"). Tenéis el directorio `devicequery` comprimido en Poliformat

-Descomprimir paquete en vuestro home y ejecutar el makefile

Averiguando características de nuestra(s) GPU(s)

Devicequery muestra por pantalla información sobre la(s) GPUs disponibles

CUDA DeviceQuery (Runtime API) versión (CUDART static linking)

Detected 1 CUDA Capable device(s)

Device0: "Tesla K20c"

CUDA Driver Version/ RuntimeVersion7.0 / 7.0

CUDA Capability Major/Minor versión number: 3.5

Total amount of global memory: 4800 MBytes(5032706048 bytes)

(13) Multiprocessors, (192) CUDA Cores/MP: 2496 CUDA Cores

GPU Max Clock rate: 706 MHz (0.71 GHz)

Memory Clock rate: 2600 Mhz

Memory Bus Width: 320-bit

L2 Cache Size: 1310720 bytes

Maximum Texture Dimension Size(x,y,z) 1D=(65536), 2D=(65536, 65536), 3D=(4096, 4096, 4096)

Maximum Layered1D Texture Size, (num) layers1D=(16384), 2048 layers

Maximum Layered2D Texture Size, (num) layers2D=(16384, 16384), 2048 layers

Total amount of constant memory: 65536 bytes Total amount of shared

memory per block: 49152 bytes

...

Programas CUDA

Aunque pueden hacerse programas en C, C++ o Fortran que acceden a las librerías, lo normal es hacer los programas con extensión .cu

Se compila con el compilador nvcc, que hace uso del compilador gcc (en Linux) o del compilador de Microsoft cl (en Windows).

Hay un buen número de opciones de compilación específicas de CUDA, aunque las típicas suelen hacer la misma tarea (-g para depurar, -O3 para optimizar, -o para darle nombre al ejecutable)

Primer programa en CUDA: Hola mundo

```
#include <stdio.h>
__global__ void kernel()
{
}
```

```
int main()
{
    kernel<<<1,1>>>();
    printf("Hello world");
    return 0;
}
```

← `__global__` indica código que se ejecuta en el "device" (GPU)

← El main se ejecuta en la CPU; el kernel se ejecuta en la GPU

En este caso, en la GPU no se hace nada. Se compila con el compilador de CUDA nvcc (que llama a gcc)

Estructura general (normal) de programa en CUDA:

Los espacios de memoria de la CPU (HOST) y de la GPU (DEVICE) son diferentes; para poder usar la GPU necesitamos funciones para reservar memoria (CudaMalloc), liberar memoria (CudaFree), Copiar memoria de CPU a GPU y de GPU a CPU (CudaMemcpy)

Programa CPU
(main)

Inicializar memoria CPU
(malloc, leer ficheros,...)

Reservar memoria GPU
(CudaMalloc)

Enviar datos de CPU a GPU
(CudaMemcpy)

Llamar a un kernel(se ejecuta en la GPU)
(nombre_de_kernel<<<x,y>>>(...parametros))

Copiar resultados de la GPU a CPU
(CudaMemcpy)

Liberar memoria
(CudaFree)

Sumar dos enteros en CUDA

```
#include<stdio.h>
__global__ void suma(int a, int b, int *c)
{
    *c=a+b;
}
int main()
{
    int c;
    int *dev_c;

    cudaMalloc ( (void**)&dev_c, sizeof(int) );

    suma<<<1,1>>>(2,7,dev_c);

    cudaMemcpy( &c ,dev_c,  sizeof(int), cudaMemcpyDeviceToHost);

    printf("2+7 = %d\n", c);

    cudaFree(dev_c);
    return 0;
}
```

Sumar dos enteros en CUDA

Observaciones:

- 1) CudaMalloc y CudaFree equivalen a malloc y free, pero para memoria de la GPU.
- 2) Los punteros a los que se le ha dado memoria con CudaMalloc no se pueden usar "correctamente" en el código de la CPU (solo a través de llamadas a kernels y cudamemcpy).
- 3) De forma similar, un puntero al que se le da memoria con malloc no se debe usar en la GPU.
- 4) Podemos copiar memoria (datos) entre CPU y GPU usando CudaMemcpy con las opciones cudaMemcpyDeviceToHost o cudaMemcpyHostToDevice.
- 5) No hace falta reservar memoria para parámetros de tipo simple (no vectores) que se pasen por valor (a,b en el ejemplo anterior). Si los argumentos son vectores, sí que hay que reservar memoria