

Simulación del modelo de Potts de q estados en GPUs

Juan P. de Francesco¹, Ezequiel E. Ferrero^{1,2,3}, Nicolás Wolovick¹, Sergio A. Cannas^{1,2}

¹ Facultad de Matemática, Astronomía y Física; Universidad Nacional de Córdoba, Ciudad Universitaria, 5000 Córdoba, Argentina

² Instituto de Física Enrique Gaviola (IFEG-CONICET), Ciudad Universitaria, 5000 Córdoba, Argentina

³ ferrero@famaf.unc.edu.ar

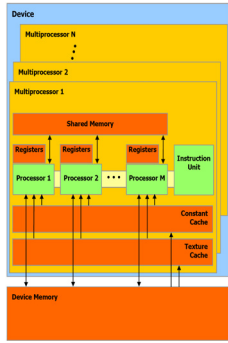
Diseñadas y desarrolladas para la industria de los video juegos, las placas de video o **GPUs** (por **Graphics Processing Units**), poseen un poder de cálculo que en la actualidad excede olgadamente al de una CPU. Esta capacidad resulta de la relativa simplicidad de la arquitectura de la GPU, comparada con la CPU, combinada por un gran número de unidades de procesamiento en paralelo en un mismo chip.

En los últimos años ha cobrado notoriedad el uso de placas de video en diversas aplicaciones científicas. Para aprovechar el potencial de las GPUs para cómputos en general es necesario programar el problema teniendo en cuenta el paralelismo inherente de las placas y la estructura jerárquica de accesos a memoria. **nVidia**, uno de los principales fabricantes de GPUs, ha desarrollado **CUDA** (siglas de **Compute Unified Device Architecture**): un compilador y un conjunto de herramientas de desarrollo que permiten a los programadores usar una variación del lenguaje de programación C para codificar algoritmos en GPUs de nVidia.

Presentamos una implementación en CUDA de simulaciones de Monte Carlo del modelo ferromagnético de Potts de q estados bidimensional. Utilizamos como estrategia de paralelización la técnica del "tablero de ajedrez". Validamos el código reobteniendo curvas de equilibrio conocidas y ciclos de histeresis asociados a la transición de 1^{er} orden para q>4. A grandes razones obtenemos una performance mayor a **150x** en el tiempo de ejecución de los programas. Discutimos en particular un problema abierto [1] que esta nueva potencia de cálculo nos permite atacar en tiempos razonables.

Arquitectura GPU Nvidia

- Tarjeta gráfica (**GPU**) compuesta por un arreglo escalable de **multiprocesadores**.
- Cada multiprocesador contiene: **8 cores** de procesadores escalares, una unidad de instrucción multithread, **memoria compartida** y un core de precisión doble.
- Cada multiprocesador contiene un **registro de 32bits** por procesador y todos los procesadores en un mismo multiprocesador comparten la shared memory.
- También existen disponibles **memorias caché** de dos tipos: **constante** y de **textura**.
- La memoria global de la GPU es compartida por todos los multiprocesadores. Es aproximadamente 10 veces más rápida que la memoria RAM de una máquina actual.

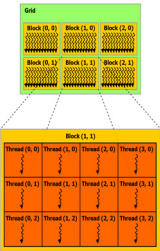


Device = GPU
Host = CPU

La GPU destina más transistores al procesamiento de datos

Programación con CUDA

- Programación en código **C** (o **Fortran**) standard usando extensiones **CUDA**.
- Cuando un programa C (o Fortran) con extensiones CUDA que corriendo en el **host** (la CPU) invoca un **kernel GPU** (sinónimo de una función GPU), un largo número de copias de este kernel, **"threads"** o **"hilos"**, son distribuidas en los multiprocesadores del **device** (la placa), en donde se ejecutan.
- Una **grilla** de kernels se subdivide en **bloques** y cada bloque en varios hilos. Esta capa de abstracción facilita la programación.
- Una **SIMT** (single-instruction multiple thread unit) se encarga de manejar y ejecutar un gran número de threads en los cores disponibles.
- La SIMT crea, planifica y ejecuta los threads en grupos de 32.
- Un grupo de 32 threads forma un warp. Threads de un mismo warp corren en el mismo multiprocesador.

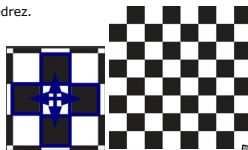


"UFO" (nuestra máquina):

Intel DX38 motherboard con Intel Core 2 Duo E8400, 2GB DDR3 1333MHz, con dos Nvidia GeForce GTX280 (240 cores y 1GB de memoria global, cada una).

Estrategia principal de paralelización (checkerboard)

- Nuestro problema es 2D y con intersecciones sólo a primeros vecinos (ver cuadro Modelo de Potts)
- La paralelización es trivial identificando a la matriz como un tablero de ajedrez.
- Todas los spins en casilleros blancos (negros) pueden ser actualizados simultáneamente (en paralelo) sin violar la secuencia temporal de la dinámica de Monte Carlo (i.e., sin violar balance detallado).
- Un paso de Monte Carlo consiste entonces en actualizar todas las blancas en paralelo e inmediatamente después todas las negras en paralelo.
- Los resultados son idénticos a los obtenidos en el Monte Carlo usual.



Características del código

- Asociamos cada spin a un hilo de ejecución diferente: **spin <-> thread**
- Utilizamos un RNG independiente para cada thread: contamos con **512x512 RNG independientes** disponibles (pares buenos de semillas).
- Subdividimos el sistema LxL en "frames" de 512x512 que mandamos al device secuencialmente.
- Para el cálculo de cantidades físicas tales como energía y magnetización utilizamos primero un cálculo local en paralelo y luego sumas recursivas en paralelo (ambas sumas tipo mariposa), acelerando el programa.

```

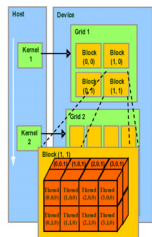
int main(int argc, char* argv[]) {
    // ...
    cudaDeviceProp prop;
    cudaGetDeviceProperties(&prop, 0);
    // ...
}

```

Pinta del código

Llamada a kernel

- Llegado el momento el programa corriendo en el host hace el llamado a un kernel que se ejecuta en el device.
- Esta función se identifica por poseer las marcas "<<<<" ">>>>", conteniendo a la estructura de grillas y bloques de threads elegida para la ejecución del kernel.
- El transpaso de información (memoria) entre el host y el device se hace explícitamente, con una directiva CUDA, sólo cuando es necesario. Esta es la operación más lenta en términos relativos.

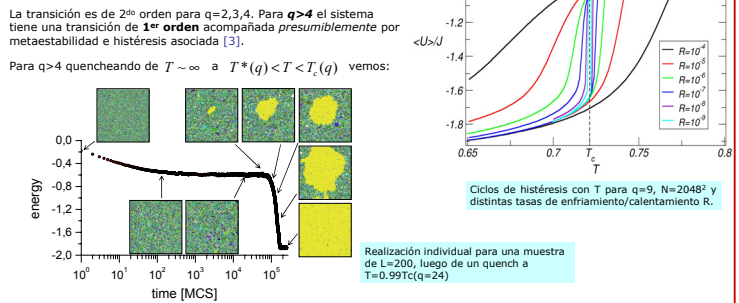


El modelo ferromagnético de Potts y el problema de Binder

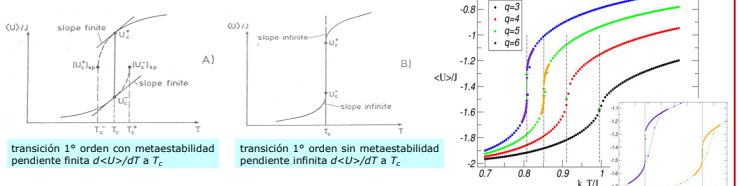
Modelo de Potts de q estados

$$\mathcal{H}_P = -J \sum_{\langle i,j \rangle} \delta(\sigma_i, \sigma_j) \quad \sigma_i = 1, 2, \dots, q \quad k_B T_c / J = 1 / \ln(1 + \sqrt{q})$$

temperatura de transición (exacta) [2]



En 1981 Binder [1] estudia (con Monte Carlo) propiedades de equilibrio para q=3,4,5,6 y propone dos escenarios posibles para la transición de 1^{er} orden (q>4) (en el límite termodinámico):

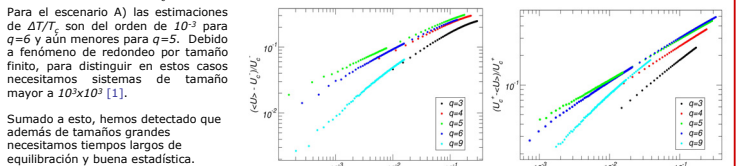


cada uno de estos escenarios serían consistentes con los comportamientos:

$$\langle U \rangle_{T < T_c} = (U_c^-)_{sp} - A^- (1 - T/T_c)^{\tau^-}$$

$$\langle U \rangle_{T > T_c} = (U_c^+)_{sp} + A^+ (1 - T_c/T)^{\tau^+}$$

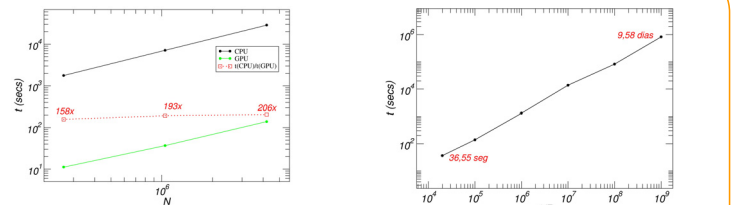
Para determinar si se cumple A) o B) la propuesta es corroborar uno de estos dos comportamientos. Dado que U_c^- y U_c^+ son conocidos exactamente [2,4], podemos medir las diferencias de energía $\langle U \rangle - U_c^-$ y $\langle U \rangle - U_c^+$ y representarlas en un gráfico log-log. Si el escenario correcto es A) deberíamos observar un crossover a una pendiente (exponente) 1 al acercarnos lo suficiente a T_c .



Sumado a esto, hemos detectado que además de tamaños grandes necesitamos tiempos largos de equilibración y buena estadística.

Creemos que el uso de las GPU nos permite ahora alcanzar el nivel de precisión necesario para atacar esta pregunta en tiempos razonables. Este es el trabajo en curso....

Performance



Tiempos de reloj característicos para un caso testigo (ciclo completo de histeresis para una muestra con q=9 entre T=0.65 y 0.8 a una tasa R=10⁻⁵) y distintos tamaños de sistema. En rojo se muestra la performance del código CUDA.

La performance de nuestro algoritmo ~200x resulta bastante superior a la obtenida recientemente para el modelo de Ising en GPUs (~60x) [5].

Referencias:

- [1] K. Binder, *J. Stat. Phys.* **24**, 69 (1981).
- [2] T. Kihara, Y. Mizuno, T. Shizume, *J. Phys. Soc. Japan* **9**, 681 (1954).
- [3] E.S. Loscar, E.E. Ferrero, T.S. Grigera, S.A. Cannas, *J. Chem. Phys.* **131**, 024120 (2009).
- [4] R.J. Baxter, *J. Phys. C* **6**, L445 (1973).
- [5] T. Preis, P. Virmann, W. Paul, J.J. Schneider, *Journal of Computational Physics* **228**, 4468 (2009)