



Máster Universitario en Computación en la Nube y de Altas Prestaciones
PROGRAMACIÓN DE GPUS CON CUDA Y OPENCL (PGPU)

Tarea 2

Introducción

Para compilar los ejercicios se puede utilizar el siguiente comando:

```
nvcc -lineinfo -Xptxas=-v -arch=sm_60 -o ejercicio ejercicio.cu
```

Ejercicio 1

En este ejercicio vamos a construir una infraestructura que vamos a utilizar en lo sucesivo para el resto de ejercicios. El código facilitado en el fichero **SumaVectores.cu** contiene dicho esqueleto. Los siguientes pasos van a consistir en el rellenado de dicho esqueleto.

1. La función **main** está completamente implementada. Se llama a la función **vector_sum**, que suma los dos vectores de tamaño **n** pasados como segundo y tercer argumentos, respectivamente, y devuelve el vector sumado como último argumento. Esta función servirá para comprobar que el resultado es correcto. La implementación es trivial y ya se encuentra implementada.
2. La función **cu_vector_sum** debe hacer lo propio en la GPU. En esta función se distinguen los punteros a la memoria de la CPU (precedidos por **h_**) de aquellos que apuntan a la memoria de la GPU (precedidos por **d_**). Véase que hay declarados tres punteros a GPU: **d_a**, **d_b** y **d_c**. Lo primero que hay que hacer es reservar memoria, apuntada por dichos punteros, en GPU mediante la función **cudaMalloc** en el lugar indicado.
3. Copiar los vectores de CPU **h_a** y **h_b** en GPU, es decir, en los vectores **d_a** y **d_b**, respectivamente, mediante la función **cudaMemcpy**.
4. Calcular los bloques de threads totales que se necesitan para “cubrir” los **n** elementos a sumar de los vectores sobre la variable **nblocks**. Por ejemplo, si **n=1300** y **blocksize=32**, el número total de bloques será de **nblocks=41**.
5. A continuación se crean dos variables, **dimGrid** y **dimBlock**, de tipo **dim3** para representar la malla de bloques y el tamaño de bloque de threads. En caso de que la malla de bloques (o del tamaño de bloque) sean “lineales” o 1D no es necesario crear esta variable pero, si lo hacemos así, nos servirá para el futuro.
6. En este paso se debe llamar al kernel mediante la notación

```
nombre_de_kernel<<<...>>>( argumentos ).
```

Obsérvese que al principio del fichero puede encontrarse la implementación inacabada de un kernel (`compute_kernel`). Se debe implementar dicho kernel utilizando los argumentos especificados. Para ello, hay que tener en cuenta que cada thread debe encargarse de sumar dos elementos de los vectores `d_a` y `d_b`, respectivamente, en un elemento del vector `d_c`. Para ello, hay que tener en cuenta que cada thread puede encargarse del elemento de índice

```
indice = threadIdx.x + blockDim.x * blockIdx.x.
```

Algo importante a tener en cuenta es que los threads “cubren” un espacio mayor que la memoria reservada para los vectores (`n`), por lo que se deberá evitar que ningún thread acceda a posiciones de memoria que no se han reservado previamente.

7. Ahora, de vuelta a la función `cu_vector_sum` y después de la llamada al kernel realizamos una copia del vector `d_c` en GPU a la CPU, o sea, en el vector `h_c` mediante la función `cudaMalloc`.
8. Terminamos la implementación de la función `cu_vector_sum` liberando la memoria creada en GPU mediante la función `cudaFree`. Este paso de liberar memoria es importante realizarlo dado que puede servir para detectar errores que podrían pasar desapercibidos.

Ejercicio 2

En este ejercicio se va a realizar la suma de dos matrices. En realidad, se trata de una generalización de la suma de vectores anterior solo que a dos dimensiones, aunque se trata más de una cuestión conceptual. La parte importante y diferenciada es que ahora vamos a trabajar con bloques de threads bidimensionales y mallas de bloques también bidimensionales.

Lo primero que vamos a tratar es de la representación de los datos. Las matrices bidimensionales en C pueden declararse con dos indirecciones de manera que el acceso al elemento i, j de la matriz A se realizaría así: `A[i][j]`. Sin embargo, nosotros vamos a utilizar otra manera. Las matrices se almacenarán en un array unidimensional. Sea pues una matriz matemática $A \in \mathbb{R}^{m \times n}$, la declaración de la misma en C la realizamos de la siguiente manera:

```
float *A = (float*) malloc ( m*n*sizeof(float) );
```

Para nosotros, las filas de la matriz A se almacenarán consecutivamente en memoria, es decir, el elemento $A_{i,j+1}$ se encuentra a continuación del elemento $A_{i,j}$. Para más sencillez, utilizaremos la siguiente notación para acceder al elemento $A_{i,j}$, por ejemplo, para asignarle un valor:

```
A( i, j ) = 4.1;
```

Evidentemente, lo anterior no es notación C. Para que funcione, es necesario que estén definidas las macros correspondientes, en este caso:

```
#define A(i,j)          A[ (j) + ((i)*(n)) ]
```

El fichero `SimpleMatrixSum.cu` contiene las macros necesarias ya definidas. Hay que prestar atención a esta definición, es decir, si quisiéramos que las matrices estuviesen almacenadas “por columnas”, por ejemplo, porque queremos utilizar las bibliotecas BLAS/LAPACK, la macro tendría el siguiente aspecto:

```
#define A(i,j)          A[ (i) + ((j)*(m)) ]
```

Esta manera de trabajar es cómoda pero tiene inconvenientes ya que debe haber coherencia entre la longitud de las filas (columnas) de la matriz declarada y de la definición de su macro correspondiente, lo que suele generar errores. También es necesario declarar una macro por matriz.

A continuación realizaremos los siguientes pasos sobre el fichero anterior:

1. La función principal se encuentra implementada. Pasamos a la función `cu_matrix_sum` que contiene la construcción de la infraestructura necesaria para llamar al kernel que resolverá el problema: `compute_kernel`. Implementamos las partes indicadas: reserva de memoria, transferencia de datos, cálculo de la malla de bloques y llamada al kernel. Es importante observar que la malla de bloques de threads debe “cubrir” completamente el espacio de $m \times n$ elementos que ocupan las matrices a sumar.
2. Implementación del kernel. En este punto hay que tener en cuenta que un thread tiene dos coordenadas dentro del bloque bidimensional de threads (`threadIdx.x` y `threadIdx.y`) y, además, cada bloque de threads tiene sus coordenadas dentro de la malla bidimensional de bloques de threads (`blockIdx.x` y `blockIdx.y`). Teniendo en cuenta que las dimensiones de los bloques son `blockDim.x` y `blockDim.y` para las dos dimensiones, respectivamente, podemos deducir fácilmente que cada thread tiene acceso a un elemento de la matriz `d_A` (o a un elemento de las otras dos, ya que todas se han almacenado de la misma manera en la GPU) utilizando, por ejemplo, las siguientes coordenadas:

```
int i = threadIdx.x + blockDim.x * blockIdx.x;
int j = threadIdx.y + blockDim.y * blockIdx.y;
```

Una vez implementado y probado el programa, hay que pasar al análisis correspondiente para ver si el alineamiento en memoria que hemos utilizado es adecuado. Este análisis se puede realizar con `computeprof` o con `nvprof`. Tomamos nota del tiempo que tarda en ejecutarse el kernel.

A continuación realizamos otra implementación del mismo kernel, esta vez ocupándonos de que los threads que tienen posiciones consecutivas en la dirección `x` accedan a posiciones consecutivas en memoria. Por ejemplo, el thread cuyo valor `threadIdx.x=17` debe acceder a la posición siguiente en memoria a la que accede el thread `threadIdx.x=16`. Una vez realizado esto se debe volver a analizar la aplicación con el profiler para ver si el tiempo de ejecución del kernel ha cambiado.

Ejercicio 3

El siguiente ejercicio sirve de repaso. Se parece mucho al primer ejercicio. La idea aquí consiste en implementar la operación `saxpy`. Esta operación viene definida en la biblioteca BLAS y se define de la siguiente manera. Dados dos vectores $x, y \in \mathbb{R}^n$ y un escalar $a \in \mathbb{R}$ la operación `saxpy` tiene la forma

$$y = a \cdot x + y.$$

El fichero `cu_saxpy.cu` contiene el esqueleto necesario para realizar la implementación de dicha operación en GPU. Dado que se trata de dos vectores lineales (1D) no va a ser necesario crear una malla bidimensional de bloques rectangulares de threads, será suficiente con crear bloques de threads 1D y una malla 1D de bloques de threads.

Ejercicio 4

Este ejercicio resuelve el mismo problema que el ejercicio anterior. En este caso, vamos a tratar la situación en la que la malla de threads no “cubre” todo el tamaño de los vectores. El código facilitado (`cu_saxpy1.cu`) pide tanto el tamaño de bloque de threads como el tamaño de malla o número de bloques. El número total de threads (`blockDim.x*gridDim.x`) puede no ser suficiente para que cada thread pueda tener asignado una posición de los vectores, tal como sucedía antes. La solución consiste en que cada thread se encargará de acceder a las posiciones:

```
threadIdx.x + blockIdx.x + blockIdx.x * blockDim.x + blockDim.x * gridDim.x
threadIdx.x + blockIdx.x + blockIdx.x * blockDim.x + 2*blockDim.x * gridDim.x
threadIdx.x + blockIdx.x + blockIdx.x * blockDim.x + 3*blockDim.x * gridDim.x
...
```

de los vectores para realizar la operación *saxpy* con esas componentes. Esto significa también que el kernel tendrá un bucle. La parte interesante de este ejercicio es averiguar la forma que tendrá dicho bucle.

Ejercicio 5

El siguiente ejercicio también es sencillo aunque no tanto como parece. La idea aquí consiste en implementar la operación *dot product* o producto escalar. Esta operación viene definida en la biblioteca BLAS y se define de la siguiente manera. Dados dos vectores $x, y \in \mathbb{R}^n$ la operación *dot* tiene la forma

$$a = \sum_{i=0}^{n-1} x_i \cdot y_i.$$

donde a es un escalar, $a \in \mathbb{R}$.

Hasta ahora el tamaño de la salida (cantidad de datos del resultado) ha sido siempre el mismo que el de la entrada. Eso ha facilitado bastante las cosas a la hora de paralelizar un algoritmo. Sin embargo, ahora nos enfrentamos al hecho de que el resultado tiene tamaño 1, es decir, es un escalar, mientras que la entrada es de tamaño n . Esto, que en paralelismo se le conoce con el nombre de *reducción*, supone un problema tanto de implementación como de eficiencia.

En un principio vamos a adoptar una solución “para salir del paso” utilizando las herramientas de las que disponemos. La solución pasará por implementar los dos kernels siguientes que podemos encontrar en el código facilitado (`cu_dot.cu`):

1. **compute_kernel1**: Este kernel, cuya signatura se puede ver en el fichero facilitado, recibe el tamaño de los vectores, los vectores a multiplicar y un vector de 32 elementos. El kernel será llamado de manera fija por una malla formada por un solo bloque de threads. Este bloque será 1D de 32 threads. La forma más fácil de implementar la solución de este kernel es aquella en la que cada thread se va a encargar de realizar la siguiente operación:

$$a_t = x_t \cdot y_t + x_{32 \cdot 1 + t} \cdot y_{32 \cdot 1 + t} + x_{32 \cdot 2 + t} \cdot y_{32 \cdot 2 + t} + \dots + x_{32 \cdot i + t} \cdot y_{32 \cdot i + t} + \dots,$$

siendo t el identificador del thread, o sea, `threadIdx.x`. Se calcularán términos mientras se cumpla $32i + t \leq n$. Cada thread guardará el resultado calculado, a_t , en una posición del vector `v` pasado como argumento, es decir, `v[threadIdx.x] = a;`.

2. **compute_kernel2**: El kernel anterior ha realizado un cálculo parcial que se encuentra almacenado en las entradas del vector `d_v`, un vector para el que se ha reservado memoria en la GPU. Esto se ha realizado así porque, al menos hasta ahora, no sabemos cómo

comunicar datos entre los threads de un bloque. Lo que sí sabemos es que los datos en memoria de la GPU son persistentes entre llamadas a kernels diferentes, es decir, mientras dura la ejecución del programa. El kernel que proponemos aquí se va a encargar de sumar los valores del vector `d_v` y devolver el resultado, que será el resultado de la operación producto escalar.

Este kernel será llamado por una malla de bloques de threads consistente en un solo bloque, igual que antes, pero ahora el bloque de threads estará formado por un solo thread. El segundo argumento del kernel será un vector de un solo elemento, que habrá sido creado antes de la llamada a dicho kernel como cualquier otro vector. El kernel se limitará a sumar los elementos del vector `v` y asignar el resultado a ese último vector de un solo elemento (`d_result`). De momento esto lo hacemos así dado que los kernels no pueden devolver un valor, siempre devuelven `void`.

Ejercicio 6

Este ejercicio es el mismo que el anterior pero ahora, la idea es racionalizar el uso de la variable `d_result` ya que, tratándose de una variable simple, no tiene sentido tratarla como un vector y reservar espacio en memoria dinámicamente. Lo primero que haremos es copiar el archivo `cu_dot.cu` en el archivo `cu_dot1.cu` para conservar ambos en caso de que se pidan.

1. La modificación consiste en declarar una variable en espacio global de la siguiente manera:

```
__device__ float d_result;
```

lo que implica que la variable anterior `d_result` debe desaparecer allí donde se había utilizado.

2. Esta nueva variable la vamos a inicializar a 0 (aunque no es necesario en este caso, pero así vemos cómo se hace). Para ello, utilizaremos la función `cudaMemcpyToSymbol`.
3. Seguidamente actualizamos la utilización de la variable dentro del kernel correspondiente. Ahora no se recibirá dicha variable como argumento ya que puede accederse como variable global que es.
4. Una vez calculado el valor de la variable `d_result` hay que devolver su valor a la CPU mediante la utilización de la función `cudaMemcpyFromSymbol`.

Ejercicio 7

Siguiendo con el ejemplo anterior, ahora vamos a copiar el fichero `cu_dot1.cu` en el `cu_dot2.cu`, y vamos a realizar las siguientes modificaciones:

1. De la misma manera que hemos creado un variable simple en la memoria del dispositivo (`d_result`) podemos también declarar un vector, en este caso, el vector `d_v` con 32 elementos. Esto se puede hacer porque se trata de un vector de tamaño constante, conocido en tiempo de compilación. Esta vez no es necesario que inicialicemos el vector `d_v` a ningún valor en particular. Modificamos consecuentemente los kernels que acceden a dicho vector, haciendo desaparecer el vector anterior. Observaremos que, ahora, el segundo kernel no tiene argumentos.

2. La segunda modificación consiste en unir los dos kernels en uno solo. Siempre va a ser más eficiente tener el número mínimo de kernels si esto es posible ya que la llamada a un kernel desde el host tiene coste. La idea es copiar el código del `compute_kernel2` en el `compute_kernel1`. Hay que tener en cuenta que el nuevo código que hemos introducido en el primer kernel solo ha de hacerlo uno de los threads, pongamos que es el `threadIdx.x=0`.

(No estaría de más mirar el cronograma de los ejercicios anteriores para ver el coste temporal.)

Ejercicio 8

Este ejercicio es trivial en comparación con los anteriores pero servirá para comenzar a trabajar con la memoria compartida. Se realizará en el fichero `cu_dot3.cu` que generaremos copiando de `cu_dot2.cu`. Se trata de “mover” la declaración del vector `d_v` como vector en la memoria del dispositivo a la memoria compartida del kernel, es decir, realizar la siguiente declaración

```
__shared__ float d_v[32];
```

dentro del kernel. Eso es todo para este ejercicio aunque hay que tener en cuenta que, cuando se escribe en la memoria `__shared__` suele ser necesario sincronizar los threads mediante una llamada a la función `__syncthreads()`, ya que un *warp* en ejecución puede acceder a dicha memoria “antes” de que esta haya sido escrita por otro *warp*. Por lo tanto, aunque aquí no es necesario dado que solo trabajamos con un *warp*, sería recomendable introducir dicha sincronización después de que los 32 threads escriban en el vector `d_v` y justo antes de que el thread 0 realice la suma de los 32 elementos de `d_v`.

Ejercicio 9

En este ejercicio seguimos trabajando con el mismo problema, el producto escalar de dos vectores, pero ahora va a ser algo más complicado ya que queremos utilizar una malla con más bloques de threads, no solo uno.

Ahora vamos a trabajar con un kernel al que llamaremos así:

```
compute_kernel<<< nblocks, BLOCKSIZE >>>( n, d_x, d_y, d_v );
```

donde `nblocks` se corresponde con la cantidad de bloques necesarios para “cubrir” el tamaño de los vectores, `n`, siendo `BLOCKSIZE` igual a 32. Este último valor es una constante definida en el programa y es usual hacerlo así puesto que la memoria *shared* que se declara ha de ser constante (no puede ser dinámica) y su tamaño suele ser una función lineal del tamaño de bloque de threads, como es el caso. Como antes, `d_x` y `d_y` serán los vectores a multiplicar, mientras que `d_v` será un vector auxiliar, creado igual que los anteriores, en memoria del dispositivo y de tamaño `nblocks`. La esencia de este vector consiste en lo siguiente: cada bloque de threads calculará un resultado parcial del producto escalar, que es un escalar, y lo guardará en la posición correspondiente de este vector. Hay que tener en cuenta que los bloques de threads no pueden comunicarse entre sí, salvo a través de la memoria del dispositivo, por esa razón necesitamos dicho vector.

Para este ejercicio contamos con un código esqueleto que podemos utilizar para rellenarlo (`cu_dot4.cu`). Sean x e y los vectores a multiplicar, s un vector de 32 elementos declarado en memoria *shared* y v un vector de tamaño `nblocks` declarado en memoria del dispositivo, para la implementación del kernel se sugiere el siguiente código:

Algoritmo 1 Productor escalar.

```
1: threadIdx ← ...                                ▷ Índice del thread dentro del bloque de threads.
2: blockIdx ← ...                                ▷ Índice del bloque dentro de la malla de bloques.
3: i ← ...                                        ▷ Índice del vector x e y que va a multiplicar este thread.
4: if i < n then
5:   s_threadIdx ← xi × yi
6:   __syncthreads()
7:   if threadIdx = 0 then
8:     Sumar los 32 elementos de s en la variable a.
9:     v_blockIdx ← a
10:    if blockIdx = 0 then
11:      Sumar los nblocks elementos de v en la variable a.
12:      d_result ← a
13:    end if
14:  end if
15: end if
```

Puede observarse en el código anterior que todos los threads realizan la operación de la línea 5. Después, solo los threads con índice 0 realizan la operación de las líneas 7–14, que consiste en sumar los valores de la memoria *shared* (línea 8) y guardar el resultado en la memoria del dispositivo (línea 9). Finalmente, solo el thread 0 del bloque 0 (líneas 10–13) suma los *nblocks* elementos del vector *v* y lo guarda en la variable global *d_result*.

Una vez implementado el código hay que probarlo varias veces y con tamaños de vector diferentes y preferiblemente grandes. Observaremos que para un mismo tamaño de vector el error que devuelve puede diferir de una ejecución a otra. También veremos valores excesivamente grandes para el valor del error. Eso indica que el programa no funciona bien. Es importante realizar una reflexión acerca de por qué esto es así antes de avanzar. De momento, lo dejamos ahí.

Ejercicio 10

En este ejercicio vamos a implementar un kernel que realice una transposición de matrices en la GPU. Dada la matriz $A \in R^{m \times n}$ el kernel obtiene la matriz $B \in R^{n \times m}$ tal que

$$B = A^T,$$

esto es, una transposición “out of place”.

El código proporcionado (`MatrixTransposition_v1.cu`) lee los valores *m* y *n*. A continuación, transpone la matriz en la CPU y almacena la transpuesta en otra. La rutina `cu.transpose` hace todo lo necesario para llamar al kernel. Antes de llamar al kernel, se forma una malla bidimensional de bloques cuadrados de threads (Código 1).

```
1 // Calculate blocks grid size
2 int blocks_row = ( m + BLOCKSIZE - 1 ) / BLOCKSIZE;
3 int blocks_col = ( n + BLOCKSIZE - 1 ) / BLOCKSIZE;
4 // Execute the kernel
5 dim3 dimGrid( blocks_col, blocks_row );
6 dim3 dimBlock( BLOCKSIZE, BLOCKSIZE );
7 compute_kernel<<< dimGrid, dimBlock >>>( m, n, d_A, d_B );
```

Código 1: Llamada al kernel de transposición de matriz.

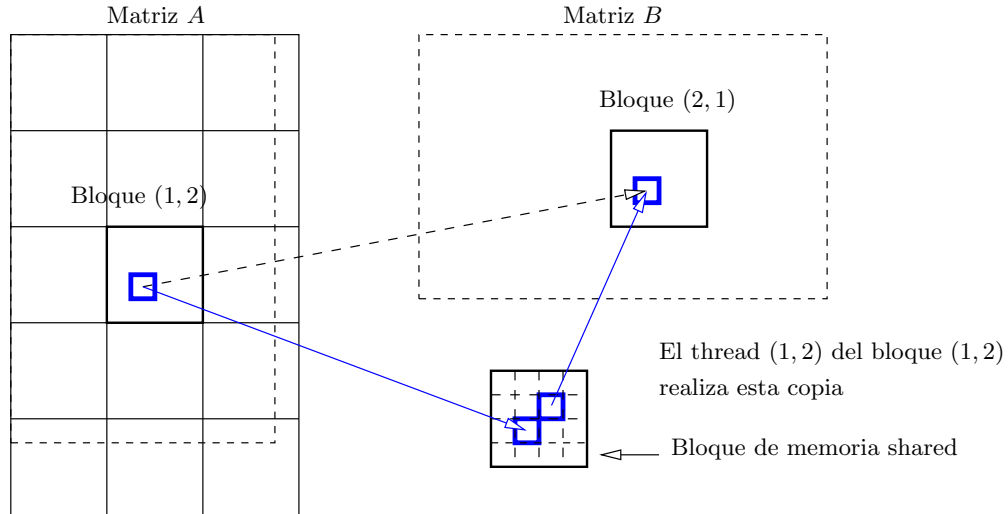


Figura 1: Ejemplo de transposición de una matriz A de 17×11 a una matriz B de 11×17 en una GPU utilizando memoria compartida. El rectángulo formado por líneas punteadas representa la matriz. Un thread (Thread(1,2)) de un bloque (Block(1,2)) copia un elemento de la matriz A ($A(11,6)$) en un elemento de la matriz B ($B(6,11)$) a través de la memoria compartida.

Para implementar el kernel sigue los pasos siguientes:

1. Obtener el índice i global de una fila de la matriz d_A .
2. Obtener el índice j global de una columna de la matriz d_A .
3. Copiar el elemento $d_A(i,j)$ en $d_B(j,i)$ para formar la matriz transpuesta. Hay que tener cuidado de no acceder a elementos más allá de los límites de la matriz.

El código que se ha implementado no accede correctamente a las posiciones de memoria de ambas matrices (puede que acceda correctamente a unas direcciones pero no a las de otra). Tomad nota del tiempo de transposición de una matriz cuadrada de tamaño 8192 para tener una referencia¹.

Ahora, vamos a utilizar una versión diferente que permita acceder a elementos consecutivos de ambas matrices. Para ello, tenemos que utilizar memoria *shared*.

La nueva versión se implementará en el fichero `MatrixTransposition.v2.cu`. Sigue los siguientes pasos:

1. Obtener las coordenadas (x,y) del thread en ambas dimensiones dentro del bloque.
2. Guardar los índices del bloque en dos variables, por ejemplo, `BLOQUE_X` y `BLOQUE_Y`.

Observar la Figura 1 para entender los siguientes pasos. En esta versión, utilizaremos un valor constante para el tamaño de bloque `BLOCKSIZE`. Esto es porque vamos a utilizar ese valor para declarar memoria compartida y su tamaño de asignación debe ser conocido en tiempo de compilación. Como se ve, para simplificar se está utilizando un bloque cuadrado. A la izquierda de la figura podemos ver la matriz A . Supongamos que un bloque de threads está procesando el bloque de matriz (1,2). Como las matrices se almacenan por filas, utilizaremos la primera dimensión del bloque de threads (`blockIdx.x`) para acceder a los elementos de fila de A . Con el

¹Por ejemplo, `nvprof ./MatrixTransposition.v1 16384 16384`.

fin de lograr accesos a memoria *coalescentes*, el mismo bloque de threads almacenará el bloque $(2, 1)$ de la matriz B . Obsérvese que una entrada de bloque de A ($A(i, j)$) y la misma entrada en B ($B(i, j)$), ambas deben ser gestionadas por el mismo thread dentro del bloque de threads para poder acceder a ambos elementos de la matriz por filas (como si se tratara de una mera copia sin transposición). Para transponer los elementos de ese bloque, utilizaremos una copia intermedia en memoria compartida.

Continúa con los pasos siguientes. Ten en cuenta que el índice de columna debe recorrer entradas consecutivas de la matriz y que la diferencia entre la matriz A y B cae en el índice de bloque de threads.

3. Obtener el índice global a una fila de A (llamarlo `r_A`).
4. Obtener el índice global de una columna de A (llamarlo `c_A`).
5. Obtener el índice global de una fila de B (llamarlo `r_B`).
6. Obtener el índice global de una columna de B (llamarlo `c_B`).
7. Declarar un bloque de memoria compartida de dimensión `BLOCKSIZE×BLOCKSIZE`.
8. Dentro de una cláusula `if` que impide acceder a los elementos de A más allá de los límites $m \times n$, escribir la copia del elemento correspondiente de A a B a través de memoria compartida. Recuerda sincronizar los threads una vez que los datos se han guardado en la memoria compartida (`__syncthreads()`).

Utiliza el profiler de CUDA (p.e. `nvprof`) de nuevo para analizar el tiempo de ejecución y comparar con la versión anterior. Ten en cuenta que, por simplicidad, vamos a utilizar siempre un tamaño de matriz que sea múltiplo de `BLOCKSIZE`.

Para terminar, haced la prueba de declarar memoria *sin conflicto de acceso a bancos de memoria*.