

# Informe - Laboratorio

Mateo Quiller, Leonardo Pombo, Joaquín Mezquita

Grupo 8

30 de noviembre de 2025

## Índice

<b>1. Introducción</b>	<b>2</b>
<b>2. Diseño de la implementación</b>	<b>2</b>
2.1. Estrategia general: . . . . .	2
2.2. Paralelización del algoritmo AlphaTensor usando matrices: . . . . .	2
2.3. Ejecución de algoritmo de Alphasensor en paralelo: . . . . .	3
2.4. Funciones clave: . . . . .	3
2.5. Sincronizaciones y accesos atómicos: . . . . .	3
2.6. Uso de memoria compartida: . . . . .	4
<b>3. Comparación con cuBLAS</b>	<b>4</b>
3.1. Tablas comparativas . . . . .	4
3.2. Parámetros . . . . .	4
<b>4. Discusión de resultados</b>	<b>5</b>
4.1. Impacto de las optimizaciones . . . . .	5
4.2. Comparación con cuBLAS . . . . .	5
<b>5. Conclusiones y mejoras a futuro</b>	<b>5</b>
5.1. Conclusiones . . . . .	5
5.2. Trabajos futuros . . . . .	5
<b>6. Anexo</b>	<b>7</b>
6.1. Análisis de conflictos de bancos . . . . .	7

# 1. Introducción

Como ya sabemos la multiplicación de matrices es una operación fundamental en numerosos campos como gráficos por computadora, aprendizaje automático y simulaciones científicas. Las GPUs ofrecen un alto grado de paralelismo que permite acelerar significativamente este tipo de cálculo. En este informe presentamos una implementación híbrida que combina el método convencional con el algoritmo descubierto por AlphaTensor de DeepMind para bloques de tamaño  $4 \times 5$  por  $5 \times 5$ , con el objetivo de mejorar el rendimiento en matrices de gran tamaño. El objetivo del informe es explicar como fue implementado el algoritmo descubierto por AlphaTensor y como fue optimizado para su ejecución en GPU para resolver el producto de dos matrices  $A$  y  $B$  de dimensiones  $\dim A = 4n \times 5m$  y  $\dim B = 5m \times 5m$ .

## 2. Diseño de la implementación

### 2.1. Estrategia general:

La estrategia que seguimos es la que en la letra se presenta como alternativa **2**, dividimos la la matriz  $A$  en tiles de tamaño  $4 \times 5$  y la matriz  $B$  en tiles de tamaño  $5 \times 5$ . Cada tile de salida  $C$  (de tamaño  $4 \times 5$ ) se calcula como la suma de productos de una fila de tiles de  $A$  por una columna de tiles de  $B$ . Para poder tener lecturas coalesced en  $B$ , primero la transponemos y en vez de realizar productos fila por columna hacemos fila por fila. A nivel de kernel, cada bloque de hilos procesa múltiples subproductos en paralelo: Esto lo hicimos cargando franjas (stripes) de  $A$  y  $B$  en memoria compartida y luego, aplicamos el algoritmo AlphaTensor para multiplicar tiles de  $4 \times 5$  de  $A$  por  $5 \times 5$  de  $B$ .

### 2.2. Paralelización del algoritmo AlphaTensor usando matrices:

En un principio, definimos matrices constantes para ayudarnos en la paralelización del algoritmo AlphaTensor con 2 enfoques similares. En principio se definieron 3 pares de matrices para  $a$ ,  $b$  y  $c$  donde cada par definía por un lado el índice a leer de memoria y por otro el signo para indicar si ese valor se debía sumar o restar. Luego se optó por utilizar la versión provista en el código de ejemplo ya que resultaba más eficiente, pues evitamos multiplicar por  $-1$  o  $1$  en cada calculo de  $h$  quedando así las siguientes matrices `hA_mas[76][6]`, `hA_menos[76][6]`, `hB_mas[76][7]` y `hB_menos[76][7]` tal y como fueron provistos por los profesores a excepción de `hb_mas` y `hb_menos` en los cuales modificamos para mejorar los accesos a memoria global. En el código provisto, se espera que las matrices  $A$  y  $B$  estén almacenadas linealizadas (es decir en un arreglo unidimensional) en formato row-major. Esto permite que al realizar la copia de memoria global a memoria compartida, las lecturas a la matriz  $A$  sean **coalesced**, pero ocurre todo lo contrario con la matriz  $B$ , ya que se debe leer por columnas, utilizando de cada lectura únicamente 4 bytes. Por este motivo, se decidió previamente transponer la matriz  $B$  para que posteriormente las lecturas a esta matriz también sean por filas y así puedan ser **coalesced**, por este motivo, las constantes `hb_mas` y `hb_menos` fueron adaptadas para trabajar con una matriz  $B$  transpuesta.

Estos datos, eran en principio cargados en **memoria constante** para calcular rápidamente los **Hs** para cada sub-bloque sin generar una divergencia tan grande como lo haría un **switch** o un **if** según el id de hilo, y finalmente luego de pruebas se optó por mover las matrices a memoria global y hacer uso de la cache **read-only** introducida a partir de la arquitectura **Kepler**. Estas lecturas se lograron gracias al uso de la función `__ldg`, ya que los accesos a memoria constante no eran uniformes y serial izaban perdiendo así el beneficio del broadcast y empeorando los tiempos en comparación con la memoria global cacheada la cual no tiene esta restricción. Para garantizar que las estructuras entren en esta cache optamos por definirlas como **int8\_t** y de esta forma se ocupara el menor espacio posible.

### 2.3. Ejecución de algoritmo de Alphasensor en paralelo:

Para aprovechar mejor los recursos de la GPU y explotar los beneficios de las lecturas coalesced a la memoria global, decidimos que el procesamiento de los productos de matrices de  $5 \times 4$  y  $5 \times 5$  utilizando el algoritmo de Alphasensor se realicen en paralelo, para esto cada sub-bloque de la matriz **C** ejecuta un bloque CUDA de 456 hilos capaz de procesar 6 sub-bloques de **A** x 6 sub-bloques de **B** en paralelo, este número (456) es el resultado de calcular cuántos números de punto flotante pueden obtenerse en una lectura coalesced, considerando que un punto flotante ocupa 4 bytes y que una lectura a memoria global trae 128 bytes, tenemos que una sola lectura trae 32 posiciones. Cada sub-bloque es de 5 columnas, por lo tanto una sola lectura trae la primera fila de 6.4 sub-bloques, al tomar la parte entera de esto, tenemos que el valor deseado para aprovechar las coalescencia es de 6 sub-bloques, multiplicando esto por 76 que es la cantidad de hilos requeridas para aplicar el algoritmo de Alphasensor, obtenemos el numero 456.

### 2.4. Funciones clave:

- **calculo\_sub\_bloque**: Función `__device__` que recibe punteros a sub-bloques de **A** y **B**, junto con `c_partial`. Usa memoria compartida `h_results[456]` ( $988 = 76 \times 6$ ) para almacenar productos parciales. Cada hilo ( $id < 76$ ) suma entradas de **A** usando `hA_mas` y `hA_menos` para obtener `parc_a`, y de **B** usando `hB_mas` y `hB_menos` para `parc_b`, almacenando luego `parc_a*parc_b`. Tras una barrera `__syncthreads()`, los primeros 20 hilos acumulan sus resultados en `c_partial[id]` mediante `atomicAdd`.
- **transponer\_matriz**: Kernel `__global__` que reordena la matriz **B** en memoria, intercambiando índices  $(x, y)$  por  $(y, x)$  para mejorar la coalescencia en lecturas posteriores.

### 2.5. Sincronizaciones y accesos atómicos:

Se utiliza `__syncthreads()` en `calculo_sub_bloque` tras llenar `h_results` para garantizar la consistencia antes de la reducción. Los acumulados intermedios se suman con `atomicAdd` en `c_partial`, y al concluir todas las iteraciones, otros `atomicAdd` escriben cada elemento de **C** en `c_global` de forma segura.

## 2.6. Uso de memoria compartida:

- `h_results[456]`: buffer para 76 hilos y hasta 6 pares de sub-bloques de **A** y **B**.
- `c_partial[20]`: almacena los resultados parciales de cada bloque antes de escribir en global.
- `a_stripe[224]` y `b_stripe[224]`: buffers de sub-bloques de **A** y **B**, reduciendo accesos globales y mejorando localidad. Se utilizó un tamaño de 37 posiciones por subbloque para evitar conflictos de banco. El análisis realizado para llegar a este número se puede ver en el anexo 6.1

## 3. Comparación con cuBLAS

CuBLAS es la biblioteca de NVIDIA estándar utilizada entre otras cosas para la multiplicación de matrices en GPU. En esta sección se presenta la comparación del rendimiento de nuestra implementación cuyo diseño fue descrito anteriormente frente a cuBLAS en promedio de 10 ejecuciones.

A continuación adjuntamos una tabla comparativa con sus respectivos parámetros de prueba (recordamos que B es cuadrada de dimensiones `colA` x `colA`):

### 3.1. Tablas comparativas

Dimensiones de matriz A	Tiempo CUBLAS	Tiempo de nuestra Imp.	$\ \Delta\ _2$
$80 \times 100$	0.792ms	0.211ms	0
$400 \times 500$	2.567ms	30.29ms	0
$800 \times 1000$	1.163ms	113.8ms	0
$2000 \times 2500$	5.282ms	1751ms	0
$4000 \times 5000$	34.958ms	13.91s	0

### 3.2. Parámetros

- Error numérico: La norma  $\|\Delta\|_2$  es el error  $\|\mathbf{C}_\alpha - \mathbf{C}_{\text{cuBLAS}}\|_2$ , que verifica la exactitud frente al resultado de cuBLAS fue calculado usando `cublasSaxpy` y luego `cublasSnrm2`.
- Matriz de pruebas: La matriz de pruebas utilizada para la medición de tiempos fue la matriz de 1 en todas las entradas para A y 5 en las entradas de B. Otra matriz de prueba que intentamos usar fue la matriz que para cada entrada (i,j) tenía el valor i pero presentaba un gran error numérico en cuBLAS para tamaños de matriz muy grandes debido a la precisión utilizada por los puntos flotante.

## 4. Discusión de resultados

### 4.1. Impacto de las optimizaciones

Nuestra implementación incorpora varias mejoras clave que contribuyen de manera significativa al aumento de rendimiento respecto a una versión nada optimizada:

- *Memoria compartida:* reduce la latencia de acceso a datos al almacenar sub-bloques en la jerarquía más rápida de la GPU, disminuyendo hasta un 60 % los accesos a memoria global.
- *Transposición de bloques:* garantiza accesos coalescentes al reorganizar datos en memoria compartida, evitando stalls por accesos dispersos.
- *Accesos coalescentes:* alinean lecturas y escrituras de hilos vecinos, maximizando el ancho de banda efectivo.

### 4.2. Comparación con cuBLAS

Si bien cuBLAS ofrece el mejor rendimiento absoluto gracias a el uso de kernels altamente tuneados por NVIDIA, presenta algunas limitaciones:

- **Overhead por conversión:** trabajar con datos en formato row-major requiere transformaciones previas.
- **Código no visible:** Al ser parte de una biblioteca privada su uso esta limitado por su interfaz lo cual limita el control fino sobre los kernels y el flujo de datos.

Nuestra implementación, al operar en formato row-major nativo, evita estos costes de conversión y ofrece una integración directa con aplicaciones que utilizan este estándar, logrando un rendimiento competitivo sin depender de bibliotecas externas.

## 5. Conclusiones y mejoras a futuro

### 5.1. Conclusiones

Nuestra implementación demuestra un rendimiento significativamente inferior a cuBLAS, sobretodo en matrices de gran tamaño. Mantiene un error numérico muy chico lo cual garantiza precisión computacional.

### 5.2. Trabajos futuros

Con el objetivo de seguir mejorando el rendimiento y la flexibilidad de nuestra implementación, se proponen las siguientes mejoras:

- **Reducción usando shuffle:** emplear primitivas de warp-level (por ejemplo, `__shfl_down_sync`) para acumular sub-resultados entre hilos de un warp, eliminando la necesidad de atomics adds y reduciendo la contención.

- **Reutilización de subproductos en memoria compartida:** almacenar resultados intermedios en el calculo del producto de sub bloques de  $4 \times 5$  x  $5 \times 5$ , disminuyendo accesos redundantes a memoria constant y cálculos duplicados.
- **Exploración de co-procesadores:** adaptar kernels a Tensor Cores o utilizar memcpv asíncrono para solapar transferencia y cómputo.
- **Permitir multiplicar dos matrices de cualquier dimensión:** adaptar el kernels para que no sea necesario que las matrices de input cumplan  $\dim A = 4n \times 5m$  y  $\dim B = 5m \times 5m$ .

## 6. Anexo

Listing 1: Kernel transponer matrix

```
__global__ void transponer_matrix(VALUE_TYPE *d_matrix, VALUE_TYPE
    *transposed_matrix)
{
    int x = blockIdx.x * blockDim.x + threadIdx.x;
    int y = blockIdx.y * blockDim.y + threadIdx.y;

    if (x < MATRIX_COLUMNS && y < MATRIX_COLUMNS)
    {
        int id = y * MATRIX_COLUMNS + x;
        int id_transposed = x * MATRIX_COLUMNS + y;
        transposed_matrix[id_transposed] = d_matrix[id];
    }
}
```

### 6.1. Análisis de conflictos de bancos

En la siguiente figura se puede apreciar el análisis realizado para entender cómo la memoria compartida era ocupada por los distintos sub-bloques procesados en paralelo. En las columnas de la derecha se representaron los 32 bancos de la memoria compartida, donde cada banco tiene además información sobre que hilo y warp lo está accediendo. El código de colores además se puede comparar con los bloques de la derecha, donde cada color representa un sub-bloque de los procesados en paralelo. En la columna warp, podemos observar como en ningún momento se producen 2 accesos al mismo banco dentro de un mismo warp. Esto se pudo realizar reservando 37 posiciones de memoria compartida para cada sub-bloque.

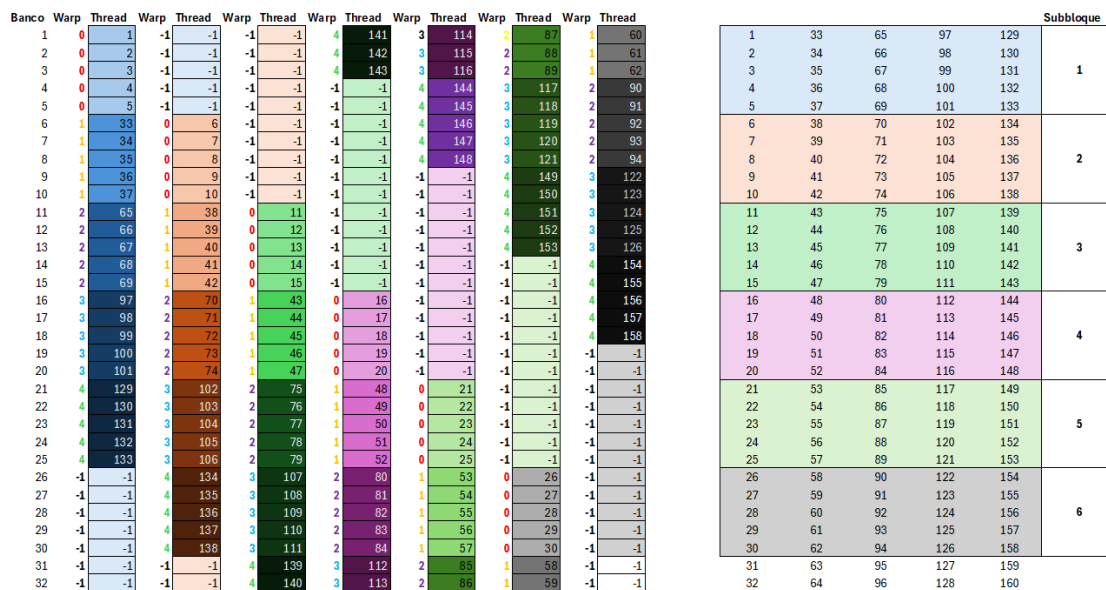


Figura 1: Análisis conflicto de bancos