

COALESCED MEMORY

OPTIMIZACIÓN DE MEMORIA EN GPUS

Jose Maureira



QUÉ ES COALESCED MEMORY?

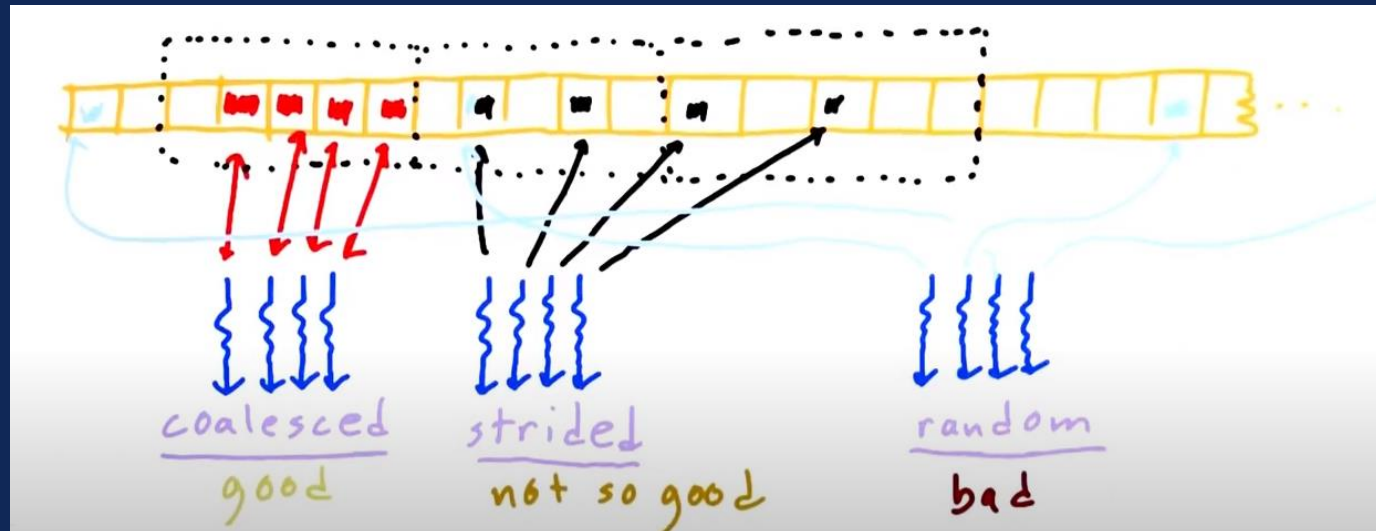
es una técnica en GPUs que maximiza la eficiencia al permitir que los hilos de un *warp* accedan a direcciones de memoria global contiguas. Esto permite al hardware agrupar los accesos en una sola transacción, reduciendo el tiempo de acceso y mejorando el rendimiento.

CÓMO FUNCIONA?

Acceso en Bloques: La GPU accede a grandes bloques de memoria, incluso si un hilo solo necesita una pequeña parte. Esto significa que si otros hilos necesitan datos en el mismo bloque, el acceso se comparte y se reutiliza, haciendo el acceso más eficiente.

TIPOS DE ACCESO A MEMORIA

- **Coalesced Access (Coalescido):** Cada hilo accede a una dirección de memoria contigua, logrando que una única transacción de memoria sea suficiente para cubrir múltiples hilos. Esto maximiza el rendimiento.
- **Strided Access (Escalonado):** Cada hilo accede a una dirección de memoria con un "salto" (stride) entre ellos, lo cual reduce la eficiencia, ya que se requieren múltiples transacciones de memoria para cubrir el acceso.
- **Random Access (Aleatorio):** Cada hilo accede a una ubicación de memoria sin ninguna relación con las demás, lo que resulta en un acceso extremadamente ineficiente, ya que cada hilo necesita su propia transacción.



ESTADO ACTUAL DE COALESCEDED MEMORY EN GPUS MODERNAS

- Coalesced Memory sigue siendo crucial en GPUs modernas para aplicaciones de alto rendimiento [1].
- GPUs recientes, como NVIDIA Ampere, han mejorado el ancho de banda de memoria y caché, pero la coalescencia sigue siendo clave [2].
- Herramientas como CUDA ayudan a optimizar estos patrones de acceso, aunque el diseño de Coalesced Memory sigue siendo fundamental en aplicaciones de gran escala [3].

¿ES TAN IMPORTANTE COALESCED MEMORY?

- Coalesced Memory optimiza el rendimiento en aplicaciones intensivas en GPU, como el aprendizaje profundo [4].
- Mantener accesos coalescidos maximiza el ancho de banda de memoria en grandes volúmenes de datos, siendo esencial para evitar cuellos de botella [5].

¿ES UN PROBLEMA RESUELTO?

- En aplicaciones con patrones de acceso regulares, la coalescencia está resuelta gracias a hardware optimizado [6].
- En accesos irregulares, como en redes neuronales complejas, sigue siendo un desafío estructurar datos para maximizar la coalescencia [7].

¿HA EMPEORADO EL PROBLEMA?

- GPUs modernas, con avanzadas cachés y memoria, reducen los efectos negativos de accesos no coalescidos [2].
- El crecimiento en complejidad de las aplicaciones, especialmente en IA, hace que los patrones de acceso varíen, requiriendo diseño específico para optimizar la coalescencia [8].

BIBLIOGRAFÍA

- NVIDIA, [CUDA C Programming Guide](#)
- NVIDIA, [NVIDIA Ampere GPU Architecture](#)
- Mittal, A. A Survey of Techniques for Improving Cache Efficiency in General-Purpose GPUs
- NVIDIA, [CUDA Optimization](#)
- Ryoo et al., Optimization Principles and Application Performance Evaluation of a Multithreaded GPU Using CUDA
- NVIDIA, [CUDA C Best Practices Guide](#)
- Volkov, V. Understanding Latency Hiding on GPUs
- Markidis et al., [NVIDIA Tensor Core Programmability, Performance & Precision](#)

COALESCED MEMORY

OPTIMIZACIÓN DE MEMORIA EN GPUS

Jose Maureira

