

## IMT2112 - Algoritmos Paralelos en Computación Científica

# Arquitecturas heterogéneas

Elwin van 't Wout

12 de noviembre de 2019



PONTIFICIA  
UNIVERSIDAD  
CATÓLICA  
DE CHILE

Facultad de Matemáticas • Escuela de Ingeniería

imc.uc.cl

# Clase previa

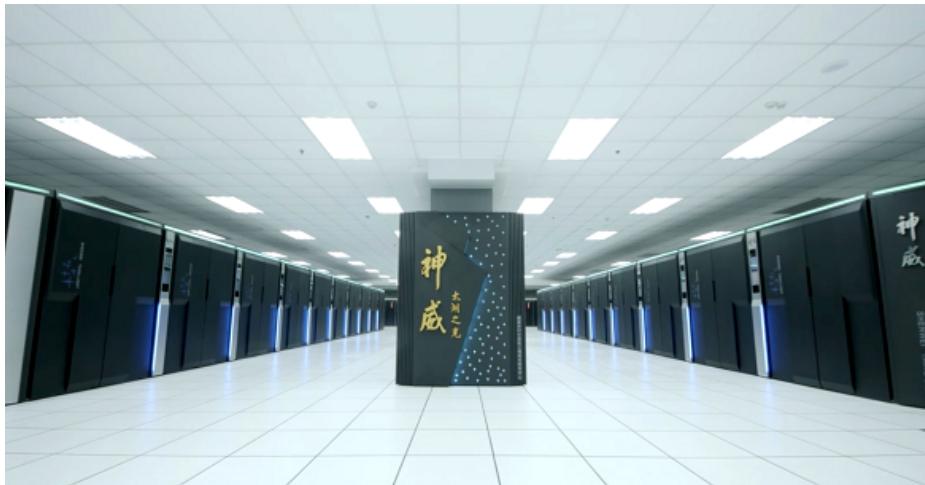
- Preacondicionamiento ILU en paralelo

# Agenda

- ¿Como usar los computadores con arquitectura heterogénea?

# Las supercomputadoras mas poderosas del mundo

La *Sunway TaihuLight* en China tiene 10.649.600 núcleos y un rendimiento de 93 petaflops



La *Summit* en los EE.UU. tiene 2.414.592 núcleos y un rendimiento de 148,6 petaflops



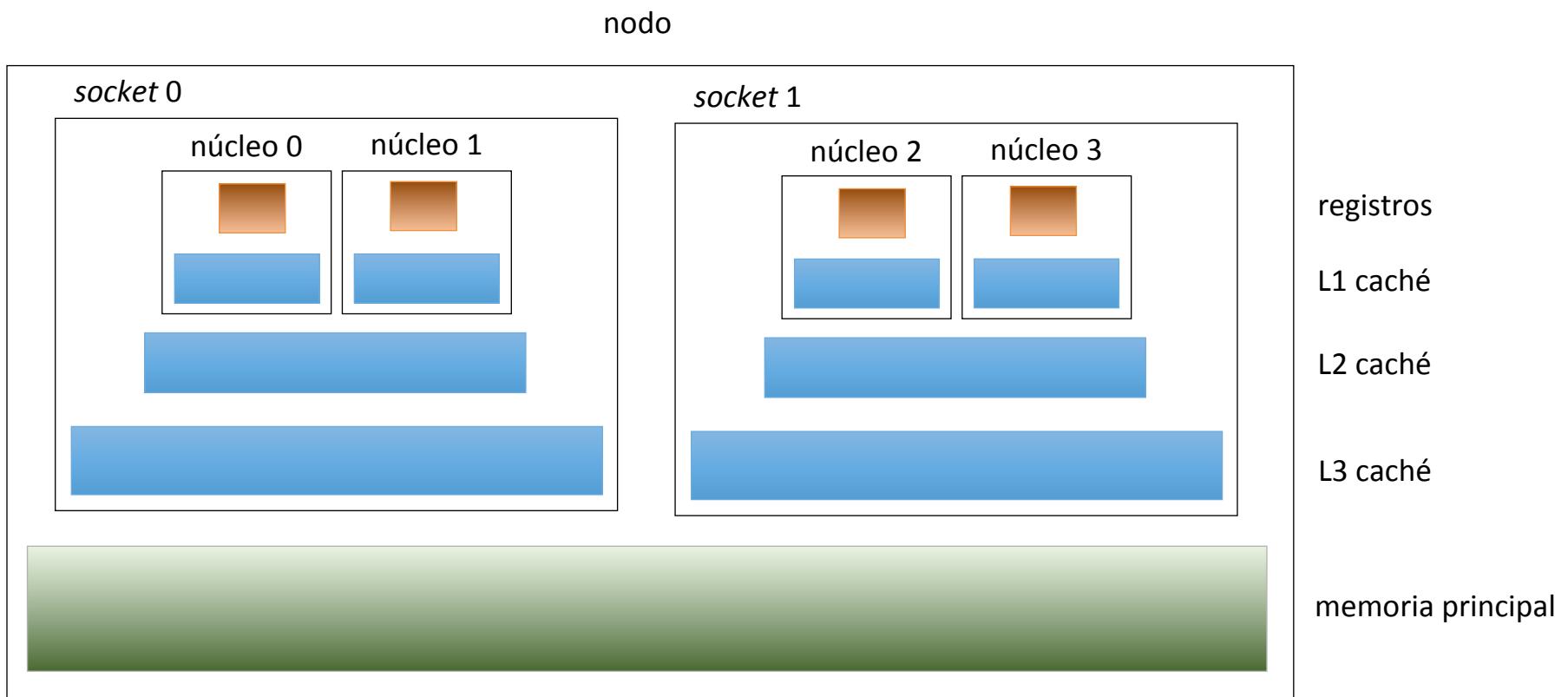
# Arquitectura de computadores

Secciones 2.3 y 2.8 del libro de Eijkhout

# La arquitectura de Von Neumann

- El primer principio es que las instrucciones se manejan como datos
- El segundo principio es manejar instrucciones como la secuencia
  - *fetch*
  - *execute*
  - *store*

# Los procesadores *multi-core*



# Procesos y hilos

- Un ‘proceso’ corresponde a la ejecución de un solo programa y contiene:
  - el código del programa (compilado)
  - el ‘*heap*’ con los datos asignados en la memoria
  - el ‘*stack*’ (‘pila’) con la información de cambio rápido:
    - contador de programa
    - elementos de datos locales
    - resultados intermedios
- Un ‘hilo’ (*thread*) es un subprocesso con su pila privada

# Procesos y hilos

- *Multiprocessing*
  - un programa que crea otros programas
  - cada programa tiene su propio espacio de datos
- *Multithreading*
  - un programa que crea varios hilos que se ejecutan simultáneamente
  - todos los hilos comparten el mismo espacio de datos
- *Hyperthreading*
  - simultáneamente ejecutando más hilos que núcleos disponibles al cambiar rápidamente entre hilos y utilizando la arquitectura de *pipelining*

# La taxonomía de Flynn

- SISD: single instruction single data
  - los procesadores tradicionales
- SIMD: single instruction multiple data
  - *vector machines* y tarjetas gráficas
- MISD: multiple instructions single data
  - no en el mercado de consumo
- MIMD: multiple instructions multiple data
  - clústers y computadores paralelos

# Arquitectura heterogénea de computadores

Sección 2.9 del libro de Eijkhout

# Sistemas heterogéneos de computación

- Los sistemas heterogéneos de computación usan diferentes tipos de procesadores simultáneamente
- Debido a limitaciones de hardware, los procesadores pueden tener diferentes características
  - el CPU es bueno para realizar instrucciones difíciles y de alto nivel, incluidas el I/O (*input* y *output*)
  - los ‘coprocesadores’, incluidas las GPU, son buenos para realizar cálculos de bajo nivel en paralelo

# Sistemas heterogéneos de computación

- Las dos arquitecturas principales para coprocesadores son GPU y MIC
- Las tarjetas gráficas (GPU - *Graphical Processing Unit*) se diseñaron para mostrar los gráficos de la computadora
- Las arquitecturas con muchos núcleos integrados (MIC - *Many Integrated Core*), como el *Intel Xeon Phi*, fueron diseñadas para realizar computación científica

# Sistemas heterogéneos de computación

- La computadora personal de IBM de primera generación (1981) utilizó dos procesadores diferentes
- NVIDIA comienza con la línea de productos Tesla GPU dirigida a la informática de alto rendimiento en 2007
  - Simultáneamente, NVIDIA introdujo la interfaz de programación CUDA
- IBM Roadrunner, la primera supercomputadora en alcanzar cálculos de petaflops en 2008
  - Usó coprocesadores Cell que también se usan en Sony Playstation 3

# Sistemas heterogéneos de computación

- Tianhe 1A (2010), la primera supercomputadora china en el lugar 1 en el TOP500
  - Utiliza tarjetas de GPU NVIDIA Tesla
- Tianhe 2 (2013), la supercomputadora más rápida del mundo durante tres años
  - Utiliza coprocesadores Intel Xeon Phi
  - El gobierno de los Estados Unidos prohíbe la venta de tecnología Intel para una actualización

# Sistemas heterogéneos de computación

- Sunway TaihuLight (2016), la supercomputadora más rápida del mundo durante dos años
  - Utiliza la arquitectura de muchos núcleos (MIC) de Sunway de construcción china
  - Chips con procesadores de 256 núcleos
- Summit (2018), actualmente la supercomputadora más rápida del mundo
  - Utiliza un sistema heterogéneo con GPU NVIDIA Volta

# Sistemas heterogéneos de computación

- En sistemas heterogéneos, todos los procesadores pueden comunicarse con la red
  - por ejemplo, CPU y GPU de uso general (*GPGPU - General Purpose GPU*)
- Sin embargo, los coprocesadores necesitan un procesador *host* para la comunicación
  - diferentes coprocesadores tampoco pueden comunicarse entre sí

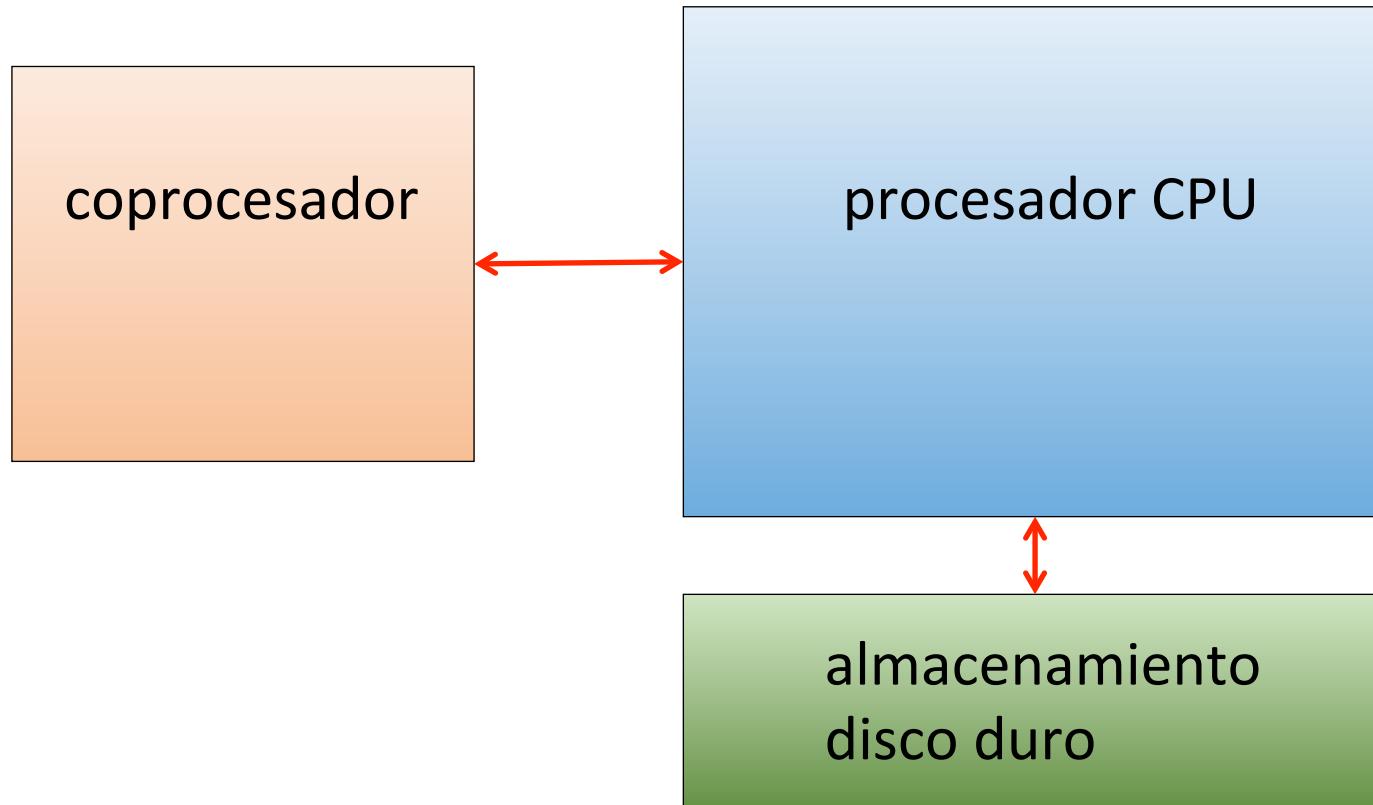
# Sistemas heterogéneos de computación

- La supercomputadora Sunway TaihuLight utiliza la arquitectura Sunway de muchos núcleos
  - cada procesador tiene 256 núcleos
  - los núcleos son de tipo RISC
- RISC - *Reduced Instruction Set Computer*
  - una solución "intermedia" entre los procesadores de GPU y CPU
  - tiene un *pipeline* mas corto que un CPU normal

# Coprocesadores

Sección 2.9 del libro de Eijkhout

# Sistemas con coprocesadores



# Sistemas heterogéneos de computación

- La idea es aprovechar las diferentes arquitecturas de procesador
  - usar CPU para realizar operaciones de alto nivel
  - utilizar los coprocesadores para descargar algoritmos de computación intensiva
- Existen diferentes arquitecturas para computadoras con coprocesadores
  - los coprocesadores comparten o no la memoria principal del procesador

# Sistemas heterogéneos de computación

- Las CPU pueden ejecutar eficientemente instrucciones avanzadas
  - manejo de entrada y salida (I/O)
  - ramas, como declaraciones *if-else*
- Las CPU tienen limitaciones para la computación científica pesada
  - Las CPU tienen un número limitado de núcleos
    - dual core: dos núcleos
    - en sistemas HPC alrededor de veinte núcleos por nodo
    - aunque el número de núcleos está aumentando
  - las CPU son lentas para cambiar entre muchos hilos
    - *hyperthreading*: dos hilos por núcleo

# Sistemas heterogéneos de computación

- El hardware de los coprocesadores está diseñado de manera que se puedan realizar muchas instrucciones al mismo tiempo
  - muchos núcleos realizan la misma instrucción en diferentes datos (SIMD)
  - cambio rápido entre hilos
- Las desventajas son:
  - comunicación lenta a la CPU a través de un *bus*
  - pequeña cantidad de memoria disponible
  - las GPU más antiguas solo admiten aritmética de precisión simple
  - generación de calor

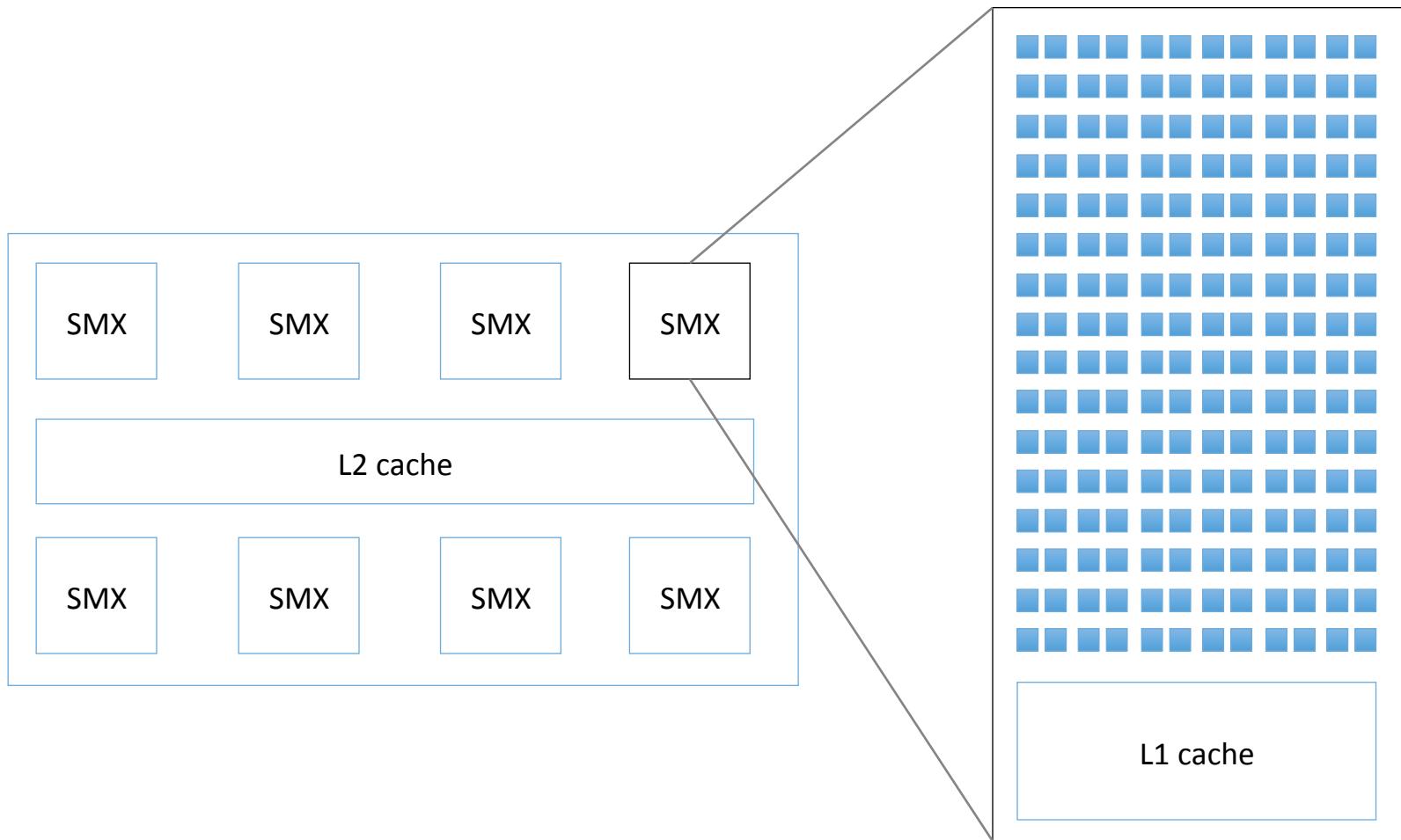
# Sistemas heterogéneos de computación

- Los procesadores de GPU NVIDIA se están volviendo ubicuos en las supercomputadoras
- En julio de 2018, Intel anunció que suspendería el desarrollo de la arquitectura MIC (Xeon Phi)

# Tarjetas GPU

Sección 2.9 del libro de Eijkhout

# Arquitectura de tarjetas GPU (NVIDIA Kepler)



# Arquitectura de tarjetas GPU

- Los componentes básicos de las GPU son multiprocesadores de transmisión (SMX - *Streaming Multiprocessor Arquitecture*)
  - procesadores con muchos núcleos, p.ej. 196 núcleos
  - caché L1 compartida, p.ej. 64 KB
  - hasta 2000 hilos
- Una tarjeta GPU tiene varias unidades SMX

# Arquitectura de tarjetas GPU

- Los hilos se ordenan en bloques de hilos
  - todos los hilos en el mismo bloque realizan la misma instrucción
  - diferentes bloques de hilos pueden realizar diferentes instrucciones
  - todos los hilos tienen datos privados y compartidos
- Combinación de modelo SIMD y SPMD
  - *Single Instruction Multiple Data*
  - *Single Program Multiple Data*
- *Single Instruction Multiple Threads* (SIMT)

# Arquitectura de tarjetas GPU

- A nivel de software, todos los bloques de hilos ejecutan la misma instrucción en diferentes datos
- A nivel de hardware, el SMX ejecuta una *warp* de 32 hilos dentro de un bloque simultáneamente

# Arquitectura de tarjetas GPU

- Hay muchos núcleos
  - Un mínimo de 32 hilos que realizan la misma instrucción simultáneamente
  - Debe tener un alto nivel de paralelismo de datos para que el paradigma SIMD sea eficiente

# Arquitectura de tarjetas GPU

- Cambio de contexto rápido
  - Cada hilo tiene su propio registro
    - Alrededor de 64 000 registros en SMX
  - Los subprocessos inactivos esperan datos y pueden iniciarse casi de inmediato cuando llegan los datos
  - Las variables declaradas como constantes están en un caché separado

# Arquitectura de tarjetas GPU

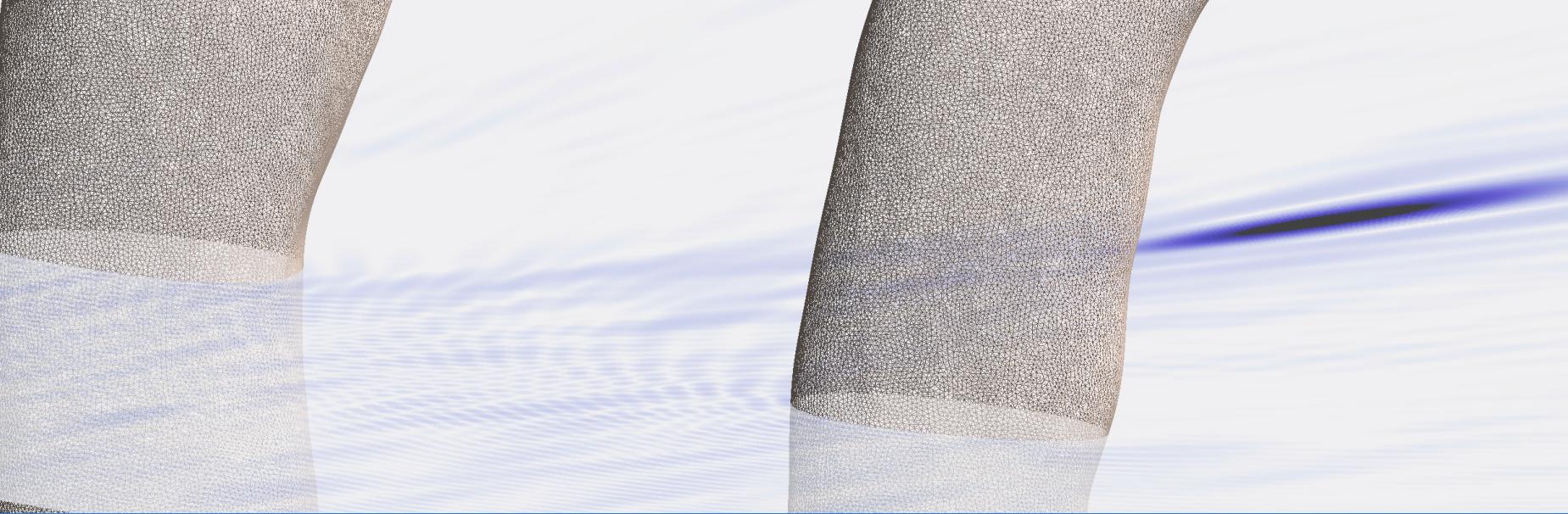
- Las GPU tienen un protocolo de coherencia relajado
  - la coherencia global no está garantizada, como en las CPU
  - no deberías tener diferentes bloques actualizando los mismos datos
  - los bloques se ejecutan en orden arbitrario

# Resumen

- Sistemas heterogéneos de computadores
- Coprocesadores
- Arquitecturas GPU
  - cada procesador tiene un gran cantidad de núcleos
  - un cambio rápido entre hilos

# Clase siguiente

- Programar en tarjetas GPU



## IMT2112 - Algoritmos Paralelos en Computación Científica

# Arquitecturas heterogéneas

Elwin van 't Wout

12 de noviembre de 2019



PONTIFICIA  
UNIVERSIDAD  
CATÓLICA  
DE CHILE

Facultad de Matemáticas • Escuela de Ingeniería

imc.uc.cl